

Another look at the “Fast Iterative Shrinkage/Thresholding Algorithm (FISTA)”

Donghwan Kim · Jeffrey A. Fessler

Date of current version: April 10, 2017

Abstract The “Fast Iterative Shrinkage/Thresholding Algorithm (FISTA)” [5], also known as a fast proximal gradient method (FPGM) in general, is widely used for efficiently minimizing composite convex functions with a nonsmooth term such as the ℓ_1 regularizer. This paper first provides an alternative way of developing FISTA (FPGM). Specifically, this paper shows that FISTA (FPGM) corresponds to an optimized approach to accelerating the proximal gradient method (PGM) with respect to the rate of decrease of the cost function. This paper then proposes a new fast algorithm called FPGM-OCG that is derived from PGM by instead optimizing the rate of decrease of the composite gradient mapping. The proof is based on the worst-case convergence analysis called Performance Estimation Problem (PEP) in [14].

1 Introduction

The “Fast Iterative Shrinkage/Thresholding Algorithm” (FISTA) [5], also known as a fast proximal gradient method (FPGM) in general, is a very widely used first-order method, since it is efficient and decreases nonsmooth composite convex cost functions with the optimal rate $O(1/N^2)$ where N denotes the number of iterations. FISTA’s speed arises from Nesterov’s accelerating technique in [25, 26] that improves the $O(1/N)$ cost function convergence rate of a proximal gradient method (PGM) [9] to the optimal $O(1/N^2)$ rate while preserving the efficient per-iteration computational complexity. Nesterov’s acceleration is interesting and effective, but has remained somewhat mysterious. Recent papers [1, 2, 7, 30, 31] have expanded our understanding of Nesterov’s acceleration.

This paper first provides an alternative way to develop Nesterov’s acceleration approach, *i.e.*, FISTA (FPGM). In particular, we show that FPGM corresponds to an optimized approach to accelerating PGM with respect to the rate of decrease of the cost function. We then propose a new fast algorithm that is derived from PGM by instead optimizing the rate of decrease of the composite gradient mapping. We call this new method FPGM-OCG (OCG for optimized over composite gradient mapping). This new method provides the best known analytical worst-case bounds for decreasing the composite gradient mapping with rate $O(1/N^{\frac{3}{2}})$ among fixed-step first-order methods. The proof is based on the worst-case convergence analysis called Performance Estimation Problem (PEP) in [14].

Drori and Teboulle’s PEP [14] casts a worst-case convergence analysis for a given optimization method and a given class of optimization problems into a meta-optimization problem. The original PEP has been intractable to solve exactly, so [14] introduced a series of tractable relaxations, focusing on first-order methods and smooth optimization problems; this PEP and its relaxations were studied for various algorithms and minimization problem classes in [15, 19, 20, 21, 22, 32, 33]. Drori and Teboulle [14] further proposed to optimize the step coefficients of a given class of optimization methods using a PEP. This approach was studied for first-order methods on unconstrained smooth convex minimization problems

This research was supported in part by NIH grant U01 EB018753.

Donghwan Kim · Jeffrey A. Fessler

Dept. of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI 48109 USA

E-mail: kimdongh@umich.edu, fessler@umich.edu

in [14, 20], and we [20] derived a new efficient first-order method, called an optimized gradient method (OGM) that has an analytic worst-case convergence bound on the cost function that is twice smaller than the previously best known bounds of [25, 26]. Recently, Drori [13] showed that the OGM exactly achieves the optimal cost function worst-case bound of first-order methods for smooth convex minimization (in high-dimensional problems).

Taylor *et al.* [32] used PEP to investigate the convergence behavior of first-order (proximal gradient) methods for minimizing nonsmooth composite convex functions. They used a tight relaxation¹ and studied the numerical convergence rates of FPGM, a proximal gradient version of OGM, and some variants versus number of iterations N . Their numerical results suggest that there exists an OGM-type acceleration of PGM that has a worst-case cost function convergence bound that is about twice smaller than that of FPGM, showing room for improvement in accelerating PGM. However, it is difficult to derive an analytical convergence bound for the tightly relaxed PEP in [32], so optimizing the step coefficients of first-order (proximal gradient) algorithms remains an open problem.

Unlike the tightly relaxed PEP in [32], this paper suggests a new (looser) relaxation of a cost function form of PEP for nonsmooth composite convex minimization that simplifies analysis and optimization of step coefficients of first-order (proximal gradient) methods, although yields loose convergence bounds. Interestingly, the resulting optimized PGM numerically appears to be the FPGM. Then, we further provide a new generalized version of FPGM using our relaxed PEP that extends our understanding of the variants of FPGM.

This paper also considers and extends the PEP analysis of the gradient norm that was discussed in [32, 33]. For unconstrained smooth convex minimization, we used such PEP to optimize the step coefficients with respect to the gradient norm in [19]. The corresponding optimized algorithm could be useful particularly when dealing with dual problems where the gradient norm convergence is important in addition to the cost function minimization (see *e.g.*, [12, 24, 28]). By extending [19], this paper optimizes the step coefficients of the first-order (proximal gradient) methods for the composite gradient mapping form of PEP for nonsmooth composite convex minimization. The resulting optimized algorithm differs somewhat from Nesterov’s acceleration and turns out to belong to the proposed generalized FPGM class.

Sec. 2 describes a nonsmooth composite convex minimization problem and first-order (proximal gradient) methods. Sec. 3 proposes a new relaxation of PEP for nonsmooth composite convex minimization problems and the proximal gradient methods, and suggests that the FPGM (FISTA) [5] is the optimized method of the cost function form of the relaxed PEP. Sec. 3 further proposes a generalized version of FPGM using the relaxed PEP. Sec. 4 studies the composite gradient mapping form of the relaxed PEP and describes a new optimized method for decreasing the norm of composite gradient mapping. Sec. 5 concludes with a comparison of the various algorithms considered.

2 Problem, methods, and contribution

We consider first-order algorithms for solving the nonsmooth composite convex minimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \{F(\mathbf{x}) := f(\mathbf{x}) + \phi(\mathbf{x})\}, \quad (\text{M})$$

under the following assumptions:

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function of the type $\mathcal{C}_L^{1,1}(\mathbb{R}^d)$, i.e., continuously differentiable with Lipschitz continuous gradient:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, \quad (1)$$

where $L > 0$ is the Lipschitz constant.

- $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is proper, closed, convex and “proximal-friendly” [9].
- The optimal set $X_*(F) = \arg \min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x})$ is nonempty, i.e., the problem (M) is solvable.

¹ Tight relaxation here denotes transforming (relaxing) an optimization problem into a solvable problem while their solutions remain the same. [32] tightly relaxes the PEP into a solvable equivalent problem under a large-dimensional condition.

We use $\mathcal{F}_L(\mathbb{R}^d)$ to denote the class of functions F that satisfy the above conditions. We additionally assume that the distance between the initial point \mathbf{x}_0 and an optimal solution $\mathbf{x}_* \in X(F)$ is bounded by $R > 0$, i.e., $\|\mathbf{x}_0 - \mathbf{x}_*\| \leq R$.

PGM is a standard first-order method for solving the problem (M) [9], particularly when the following proximal gradient update (that consists of a gradient descent step and a proximal operation [9]) is relatively simple:

$$\begin{aligned} \mathbf{p}_L(\mathbf{y}) &:= \arg \min_{\mathbf{x}} \left\{ f(\mathbf{y}) + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \phi(\mathbf{x}) \right\} \\ &= \arg \min_{\mathbf{x}} \left\{ \frac{L}{2} \left\| \mathbf{x} - \left(\mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y}) \right) \right\|^2 + \phi(\mathbf{x}) \right\}. \end{aligned} \quad (2)$$

For $\phi(\mathbf{x}) = \|\mathbf{x}\|_1$, the proximal gradient update (2) becomes a simple shrinkage/thresh-olding update, and PGM reduces to an iterative shrinkage/thresholding algorithm (ISTA) [11]. (See [9, Table 10.2] for more functions $\phi(\mathbf{x})$ that lead to simple proximal operations.)

Algorithm PGM

Input: $f \in \mathcal{F}_L(\mathbb{R}^d)$, $\mathbf{x}_0 \in \mathbb{R}^d$.

For $i = 0, \dots, N - 1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{x}_i)$$

PGM has the following bound on the cost function [5, Thm. 3.1] for any $N \geq 1$:

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \frac{LR^2}{2N}. \quad (3)$$

For simplicity in later derivations, we use the following definition of the composite gradient mapping [29]:

$$\tilde{\nabla}_L F(\mathbf{x}) := -L(\mathbf{p}_L(\mathbf{x}) - \mathbf{x}). \quad (4)$$

The composite gradient mapping reduces to the usual function gradient $\nabla f(\mathbf{x})$ when $\phi(\mathbf{x}) = 0$. We can then rewrite the PGM update in the following form reminiscent of a gradient method:

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{x}_i) = \mathbf{x}_i - \frac{1}{L} \tilde{\nabla}_L F(\mathbf{x}_i), \quad (5)$$

where each update guarantees the following monotonic cost function descent [29, Thm. 1]:

$$F(\mathbf{x}_i) - F(\mathbf{x}_{i+1}) \geq \frac{1}{2L} \|\tilde{\nabla}_L F(\mathbf{x}_i)\|^2. \quad (6)$$

For any $\mathbf{x} \in \mathbb{R}^d$, there exists a subgradient $\phi'(\mathbf{p}_L(\mathbf{x})) \in \partial\phi(\mathbf{p}_L(\mathbf{x}))$ that satisfies the following equality [5, Lemma 2.2]

$$\tilde{\nabla}_L F(\mathbf{x}) = \nabla f(\mathbf{x}) + \phi'(\mathbf{p}_L(\mathbf{x})). \quad (7)$$

This equality implies that any point $\bar{\mathbf{x}}$ with a zero composite gradient mapping ($\tilde{\nabla}_L F(\bar{\mathbf{x}}) = 0$, i.e., $\bar{\mathbf{x}} = \mathbf{p}_L(\bar{\mathbf{x}})$) satisfies $0 \in \partial F(\bar{\mathbf{x}})$ and is a minimizer of (M). As discussed, minimizing the composite gradient mapping is noteworthy for overall convergence speed in addition to decreasing the cost function. This property becomes particularly important when dealing with dual problems. In particular, it is known that the norm of the dual (sub)gradient is related to the primal feasibility (see e.g., [12, 24, 28]). Furthermore, the norm of the subgradient is upper bounded by the norm of the composite gradient mapping, i.e., for any given subgradients $\phi'(\mathbf{p}_L(\mathbf{x}))$ in (7) and $F'(\mathbf{p}_L(\mathbf{x})) := \nabla f(\mathbf{p}_L(\mathbf{x})) + \phi'(\mathbf{p}_L(\mathbf{x})) \in \partial F(\mathbf{p}_L(\mathbf{x}))$, we have

$$\begin{aligned} \|F'(\mathbf{p}_L(\mathbf{x}))\| &\leq \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{p}_L(\mathbf{x}))\| + \|\nabla f(\mathbf{x}) + \phi'(\mathbf{p}_L(\mathbf{x}))\| \\ &\leq 2L\|\mathbf{x} - \mathbf{p}_L(\mathbf{x})\| = 2\|\tilde{\nabla}_L F(\mathbf{p}_L(\mathbf{x}))\|, \end{aligned} \quad (8)$$

where the first inequality uses the triangle inequality and the second inequality uses (1) and (7). This inequality provides a close relationship between the primal feasibility and the dual composite gradient mapping. Therefore, we next analyze the convergence rate of the composite gradient mapping of PGM; Sec. 4 below discusses a first-order algorithm that is optimized with respect to the composite gradient mapping.²

The following lemma shows that PGM monotonically decreases the norm of the composite gradient mapping.

Lemma 1 *The PGM monotonically decreases the norm of composite gradient mapping, i.e., for all \mathbf{x} :*

$$\|\tilde{\nabla}_L F(\mathbf{p}_L(\mathbf{x}))\| \leq \|\tilde{\nabla}_L F(\mathbf{x})\|. \quad (9)$$

Proof The proof in [24, Lemma 2.4] can be easily extended to prove (9) using the nonexpansiveness of the proximal mapping (proximity operator) [9].

The following theorem shows the $O(1/N)$ convergence bound on the norm of composite gradient mapping for the PGM, using the idea in [28] and Lemma 1.

Theorem 1 *Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $\mathcal{F}_L(\mathbb{R}^d)$ and let $\mathbf{x}_0, \dots, \mathbf{x}_N \in \mathbb{R}^d$ be generated by PGM. Then for $N \geq 2$,*

$$\min_{i \in \{0, \dots, N\}} \|\tilde{\nabla}_L F(\mathbf{x}_i)\| = \|\tilde{\nabla}_L F(\mathbf{x}_N)\| \leq \frac{2LR}{\sqrt{(N-1)(N+2)}}. \quad (10)$$

Proof Let $m = \lfloor \frac{N}{2} \rfloor$, and we have

$$\begin{aligned} \frac{LR^2}{2m} &\stackrel{(3)}{\geq} F(\mathbf{x}_m) - F(\mathbf{x}_*) \stackrel{(6)}{\geq} F(\mathbf{x}_{N+1}) - F(\mathbf{x}_*) + \frac{1}{2L} \sum_{i=m}^N \|\tilde{\nabla}_L F(\mathbf{x}_i)\|^2 \\ &\stackrel{(9)}{\geq} \frac{N-m+1}{2L} \|\tilde{\nabla}_L F(\mathbf{x}_N)\|^2, \end{aligned}$$

which is equivalent to (10) using $m \geq \frac{N-1}{2}$ and $N-m \geq \frac{N}{2}$.

Despite its computational efficiency, PGM suffers from the slow convergence rate $O(1/N)$ for decreasing both the cost function and the norm of composite gradient mapping.³ Therefore for acceleration, this paper considers the following class of *fixed-step* first-order methods (FSFOM), where the $(i+1)$ th iteration consists of one proximal gradient evaluation, just like PGM, and a weighted summation of previous and current proximal gradient updates $\{\mathbf{x}_{k+1} - \mathbf{y}_k\}_{k=0}^i$ with step coefficients $\{h_{i+1,k}\}_{k=0}^i$.

Algorithm Class FSFOM

Input: $f \in \mathcal{F}_L(\mathbb{R}^d)$, $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{y}_0 = \mathbf{x}_0$.

For $i = 0, \dots, N-1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i) = \mathbf{y}_i - \frac{1}{L} \tilde{\nabla}_L F(\mathbf{y}_i)$$

$$\mathbf{y}_{i+1} = \mathbf{y}_i + \sum_{k=0}^i h_{i+1,k} (\mathbf{x}_{k+1} - \mathbf{y}_k) = \mathbf{y}_i - \frac{1}{L} \sum_{k=0}^i h_{i+1,k} \tilde{\nabla}_L F(\mathbf{y}_k).$$

Although the weighted summation in FSFOM seems at first to be inefficient both computationally and memory-wise, the optimized FSFOM presented in this paper have equivalent efficient recursive forms that have memory and computation requirements that are similar to PGM. Note that this class FSFOM includes PGM but excludes accelerated algorithms in [16, 27, 29] that combine the proximal operations and the gradient steps in other ways.

² One could develop a first-order algorithm that is optimized with respect to the norm of the subgradient (rather than its upper bound in Sec. 4), which we leave as future work.

³ [14, Thm. 2] and [19, Thm. 2] imply that the $O(1/N)$ rates of both the cost function bound (3) and the composite gradient mapping norm bound (10) of PGM are tight up to a constant respectively.

Among FSFOM⁴, FISTA [5], also known as FPGM, is widely used since it has computation and memory requirements that are similar to PGM yet it achieves the optimal $O(1/N^2)$ rate for decreasing the cost function, using Nesterov's acceleration technique [25, 26].

Algorithm FPGM (FISTA)

Input: $f \in \mathcal{F}_L(\mathbb{R}^d)$, $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{y}_0 = \mathbf{x}_0$, $t_0 = 1$.

For $i = 0, \dots, N - 1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i)$$

$$t_{i+1} = \frac{1 + \sqrt{1 + 4t_i^2}}{2} \quad (11)$$

$$\mathbf{y}_{i+1} = \mathbf{x}_{i+1} + \frac{t_i - 1}{t_{i+1}}(\mathbf{x}_{i+1} - \mathbf{x}_i)$$

FPGM has the following bound for the cost function [5, Thm. 4.4] for any $N \geq 1$:

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \frac{LR^2}{2t_{N-1}^2} \leq \frac{2LR^2}{(N+1)^2}, \quad (12)$$

where the somewhat mysterious parameters t_i (11) satisfy

$$t_i^2 = \sum_{l=0}^i t_l \quad \text{and} \quad t_i \geq \frac{i+2}{2}. \quad (13)$$

Sec. 3 provides a new proof of the cost function bound (12) of FPGM using a new relaxation of PEP, and illustrates that this particular acceleration of PGM results from optimizing a relaxed version of the cost function form of PEP. In addition, it is shown in [5, 8] that FPGM and its bound (12) generalize to any t_i such that $t_0 = 1$ and $t_i^2 \leq t_{i-1}^2 + t_i$ for all $i \geq 1$ with corresponding bound for any $N \geq 1$:

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \frac{LR^2}{2t_{N-1}^2}, \quad (14)$$

which includes the choice $t_i = \frac{i+a}{a}$ for any $a \geq 2$. Using our relaxed PEP, Sec. 3 further describes similar but different generalizations of FPGM that complement our understanding of FPGM.

We are often interested in the convergence speed of the norm of the (composite) gradient (mapping) in addition to that of the cost function, particularly when dealing with dual problems. To improve the rate $O(1/N)$ of the gradient norm bound of a gradient method, Nesterov [28] suggested performing his fast gradient method (FGM) [25, 26], a non-proximal version of FPGM, for the first m iterations and a gradient method for remaining $N - m$ iterations for smooth convex problems (when $\phi(\mathbf{x}) = 0$). Here we extend this idea to the nonsmooth composite convex problem (M) and use FPGM- m to denote the resulting algorithm.

Algorithm FPGM- m

Input: $f \in \mathcal{F}_L(\mathbb{R}^d)$, $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{y}_0 = \mathbf{x}_0$, $t_0 = 1$.

For $i = 0, \dots, N - 1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i)$$

$$t_{i+1} = \frac{1 + \sqrt{1 + 4t_i^2}}{2}, \quad i \leq m - 1$$

$$\mathbf{y}_{i+1} = \begin{cases} \mathbf{x}_{i+1} + \frac{t_i - 1}{t_{i+1}}(\mathbf{x}_{i+1} - \mathbf{x}_i), & i \leq m - 1, \\ \mathbf{x}_{i+1}, & \text{otherwise.} \end{cases}$$

⁴ The step coefficients of FSFOM for FPGM are [14, 20]

$$h_{i+1,k} = \begin{cases} \frac{1}{t_{i+1}} \left(t_k - \sum_{j=k+1}^i h_{j,k} \right), & k = 0, \dots, i - 1, \\ 1 + \frac{t_i - 1}{t_{i+1}}, & k = i. \end{cases}$$

The following theorem provides a $O(1/N^{\frac{3}{2}})$ convergence bound for the norm of composite gradient mapping of FPGM- m , using the idea in [28] and Lemma 1.

Theorem 2 *Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $\mathcal{F}_L(\mathbb{R}^d)$ and let $\mathbf{x}_0, \dots, \mathbf{x}_N \in \mathbb{R}^d$ be generated by FPGM- m for $1 \leq m \leq N$. Then for $N \geq 1$,*

$$\min_{i \in \{0, \dots, N\}} \|\tilde{\nabla}_L F(\mathbf{x}_i)\| \leq \|\tilde{\nabla}_L F(\mathbf{x}_N)\| \leq \frac{2LR}{(m+1)\sqrt{N-m+1}}. \quad (15)$$

Proof We have

$$\begin{aligned} \frac{2LR^2}{(m+1)^2} &\stackrel{(12)}{\geq} F(\mathbf{x}_m) - F(\mathbf{x}_*) \stackrel{(6)}{\geq} F(\mathbf{x}_{N+1}) - F(\mathbf{x}_*) + \frac{1}{2L} \sum_{i=m}^N \|\tilde{\nabla}_L F(\mathbf{x}_i)\|^2 \\ &\stackrel{(9)}{\geq} \frac{N-m+1}{2L} \|\tilde{\nabla}_L F(\mathbf{x}_N)\|^2, \end{aligned}$$

which is equivalent to (15).

As noticed by the reviewer, when $m = \lfloor \frac{2N}{3} \rfloor$, the convergence bound (15) of the composite gradient mapping roughly has its smallest constant $3\sqrt{3}$ for the rate $O(1/N^{\frac{3}{2}})$, which is better than the choice $m = \lfloor \frac{N}{2} \rfloor$ in [28].

Monteiro and Svaiter [23] considered a variant of FPGM that replaces $\mathbf{p}_L(\cdot)$ of FPGM by $\mathbf{p}_{L/\sigma^2}(\cdot)$ for $0 < \sigma < 1$; that variant, which we denote FPGM- σ , satisfies the $O(1/N^{\frac{3}{2}})$ rate for the composite gradient mapping. This FPGM- σ algorithm satisfies the following cost function and composite gradient mapping convergence bounds⁵ [23, Prop. 5.2] for $N \geq 1$,

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \frac{2LR^2}{\sigma^2 N^2}, \quad (16)$$

$$\min_{i \in \{0, \dots, N\}} \|\tilde{\nabla}_{L/\sigma^2} F(\mathbf{y}_i)\| \leq \frac{2\sqrt{3}}{\sigma} \sqrt{\frac{1+\sigma}{1-\sigma}} \frac{LR}{N^{\frac{3}{2}}}. \quad (17)$$

The convergence bound (17) of the composite gradient mapping has its smallest constant $\frac{2\sqrt{3}}{\sigma^2} \sqrt{\frac{1+\sigma}{1-\sigma}} \approx 16.2$ when $\sigma = \frac{\sqrt{17}-1}{4} \approx 0.78$, which makes the bound (17) about $\frac{16}{3\sqrt{3}} \approx 3$ -times larger than the bound (15) of FPGM- $(m = \lfloor \frac{2N}{3} \rfloor)$ at best. However, since FPGM- σ does not require one to select the number of total iterations N in advance unlike FPGM- m , the FPGM- σ algorithm could be useful in practice, as discussed further in Sec. 4.4. Ghadimi and Lan [16] also showed the $O(1/N^{\frac{3}{2}})$ rate for a composite gradient mapping convergence bound of another variant of FPGM, but the corresponding algorithm in [16] requires two proximal gradient updates per iteration, combining the proximal operations and the gradient steps in a way that differs from the class FSFOM and could be less attractive in terms of the per-iteration complexity.

FPGM has been used in dual problems [3, 4, 6, 17]; using FPGM- m and the algorithms in [16, 23] that guarantee $O(1/N^{\frac{3}{2}})$ rate for minimizing the norm of the composite gradient mapping could be potentially useful for solving dual problems. (Using F(P)GM- m for (dual) smooth convex problems was discussed in [12, 24, 28].) However, FPGM- m and the algorithms in [16, 23] are not necessarily the best possible methods with respect to the rate of decrease of the norm of the composite gradient mapping. Therefore, Sec. 4 seeks to optimize the step coefficients of FSFOM for minimizing the norm of the composite gradient mapping using a relaxed PEP.

The next section first provides a new proof of FPGM using our new relaxation on PEP, and proposes the new generalized FPGM.

⁵ The bound for $\min_{i \in \{0, \dots, N\}} \|\tilde{\nabla}_{L/\sigma^2} F(\mathbf{y}_i)\|$ of FPGM- σ is described in a big-O sense in [23, Prop. 5.2(c)], and we further computed the constant in (17) by following the derivation of [23, Prop. 5.2(c)].

3 Relaxation and optimization of the cost function form of PEP

3.1 Relaxation for the cost function form of PEP

For FSFOM with given step-size coefficients $\mathbf{h} := \{h_{i+1,k}\}$, in principle the worst-case bound on the cost function after N iterations corresponds to the solution of the following PEP problem [14]:

$$\begin{aligned} \mathcal{B}_P(\mathbf{h}, N, d, L, R) := & \max_{\substack{F \in \mathcal{F}_L(\mathbb{R}^d), \\ \mathbf{x}_0, \dots, \mathbf{x}_N, \mathbf{x}_* \in \mathbb{R}^d, \\ \mathbf{y}_0, \dots, \mathbf{y}_{N-1} \in \mathbb{R}^d, \\ \mathbf{x}_* \in X_*(F), \|\mathbf{x}_0 - \mathbf{x}_*\| \leq R}} F(\mathbf{x}_N) - F(\mathbf{x}_*) \\ \text{s.t. } & \mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i), \quad i = 0, \dots, N-1, \\ & \mathbf{y}_{i+1} = \mathbf{y}_i + \sum_{k=0}^i h_{i+1,k}(\mathbf{x}_{k+1} - \mathbf{y}_k), \quad i = 0, \dots, N-2, \end{aligned} \quad (\text{P})$$

Since (non-relaxed) PEP problems like (P) are difficult to solve due to the (infinite-dimensional) functional constraint on F , Drori and Teboulle [14] suggested (for smooth convex problems) replacing the functional constraint by a property of F related to the update such as $\mathbf{p}_L(\cdot)$ in (P). Taylor *et al.* [32, 33] discussed properties of F that can replace the functional constraint of PEP without strictly relaxing (P), and provided tight numerical convergence analysis for any given N . However, analytical solutions remain unknown for (P) and most PEP problems.

Instead, this paper proposes an alternate relaxation that is looser than that in [32, 33] but provides tractable and useful analytical results. We consider the following property of F involving the proximal gradient update $\mathbf{p}_L(\cdot)$ [5, Lemma 2.3]:

$$F(\mathbf{x}) - F(\mathbf{p}_L(\mathbf{y})) \leq \frac{L}{2} \|\mathbf{p}_L(\mathbf{y}) - \mathbf{y}\|^2 + L \langle \mathbf{y} - \mathbf{x}, \mathbf{p}_L(\mathbf{y}) - \mathbf{y} \rangle \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d \quad (18)$$

to replace the functional constraint on F . In particular, we use the following property:

$$\begin{aligned} & \frac{L}{2} \|\mathbf{p}_L(\mathbf{y}) - \mathbf{y}\|^2 - L \langle \mathbf{p}_L(\mathbf{x}) - \mathbf{x}, \mathbf{p}_L(\mathbf{y}) - \mathbf{y} \rangle \\ & \leq F(\mathbf{p}_L(\mathbf{x})) - F(\mathbf{p}_L(\mathbf{y})) + L \langle \mathbf{p}_L(\mathbf{y}) - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d \end{aligned} \quad (19)$$

that results from replacing \mathbf{x} in (18) by $\mathbf{p}_L(\mathbf{x})$. When $\phi(\mathbf{x}) = 0$, the property (19) reduces to

$$\begin{aligned} & \frac{1}{2L} \|\nabla f(\mathbf{y})\|^2 - \frac{1}{L} \langle \nabla f(\mathbf{x}), \nabla f(\mathbf{y}) \rangle \\ & \leq f\left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})\right) - f\left(\mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y})\right) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \end{aligned} \quad (20)$$

Note that the relaxation of PEP in [14, 19, 20, 21, 33] for unconstrained smooth convex minimization ($\phi(\mathbf{x}) = 0$) uses a well-known property of f in [26, Thm. 2.1.5] that differs from (20) and does not strictly relax the PEP as discussed in [33], whereas our relaxation using (19) and (20) does not guarantee a tight relaxation of (P). Finding a tight relaxation that leads to useful (or even optimal) algorithms remains an open problem for nonsmooth problems.

Similar to [14, Problem (Q')], we (strictly) relax problem (P) as follows using a set of constraint inequalities (19) at the points $(\mathbf{x}, \mathbf{y}) = (\mathbf{y}_{i-1}, \mathbf{y}_i)$ for $i = 1, \dots, N-1$ and $(\mathbf{x}, \mathbf{y}) = (\mathbf{x}_*, \mathbf{y}_i)$ for $i = 0, \dots, N-1$:

$$\begin{aligned} \mathcal{B}_{P1}(\mathbf{h}, N, d, L, R) := & \max_{\substack{\mathbf{G} \in \mathbb{R}^{N \times d}, \\ \boldsymbol{\delta} \in \mathbb{R}^N}} LR^2 \delta_{N-1} \\ \text{s.t. } & \text{Tr}\{\mathbf{G}^\top \tilde{\mathbf{A}}_{i-1,i}(\mathbf{h}) \mathbf{G}\} \leq \delta_{i-1} - \delta_i, \quad i = 1, \dots, N-1, \\ & \text{Tr}\{\mathbf{G}^\top \tilde{\mathbf{D}}_i(\mathbf{h}) \mathbf{G} + \boldsymbol{\nu} \mathbf{u}_i^\top \mathbf{G}\} \leq -\delta_i, \quad i = 0, \dots, N-1, \end{aligned} \quad (\text{P1})$$

for any given unit vector $\boldsymbol{\nu} \in \mathbb{R}^d$, by defining the $(i+1)$ th standard basis vector $\mathbf{u}_i = \mathbf{e}_{i+1} \in \mathbb{R}^N$, the matrix $\mathbf{G} = [\mathbf{g}_0, \dots, \mathbf{g}_{N-1}]^\top \in \mathbb{R}^{N \times d}$ and the vector $\boldsymbol{\delta} = [\delta_0, \dots, \delta_{N-1}]^\top \in \mathbb{R}^N$, where

$$\begin{cases} \mathbf{g}_i := -\frac{1}{\|\mathbf{y}_0 - \mathbf{x}_*\|}(\mathbf{p}_L(\mathbf{y}_i) - \mathbf{y}_i) = \frac{1}{L\|\mathbf{y}_0 - \mathbf{x}_*\|}\tilde{\nabla}_L F(\mathbf{y}_i), \\ \delta_i := \frac{1}{L\|\mathbf{y}_0 - \mathbf{x}_*\|^2}(F(\mathbf{p}_L(\mathbf{y}_i)) - F(\mathbf{x}_*)), \end{cases} \quad (21)$$

for $i = 0, \dots, N-1, *$. Note that $\mathbf{g}_* = [0, \dots, 0]^\top$, $\delta_* = 0$ and $\text{Tr}\{\mathbf{G}^\top \mathbf{u}_i \mathbf{u}_j^\top \mathbf{G}\} = \langle \mathbf{g}_i, \mathbf{g}_j \rangle$ by definition. The matrices $\tilde{\mathbf{A}}_{i-1,i}(\mathbf{h})$ and $\tilde{\mathbf{D}}_i(\mathbf{h})$ are defined as

$$\begin{cases} \tilde{\mathbf{A}}_{i-1,i}(\mathbf{h}) := \frac{1}{2}\mathbf{u}_i \mathbf{u}_i^\top - \frac{1}{2}\mathbf{u}_{i-1} \mathbf{u}_i^\top - \frac{1}{2}\mathbf{u}_i \mathbf{u}_{i-1}^\top + \frac{1}{2} \sum_{k=0}^{i-1} h_{i,k}(\mathbf{u}_i \mathbf{u}_k^\top + \mathbf{u}_k \mathbf{u}_i^\top), \\ \tilde{\mathbf{D}}_i(\mathbf{h}) := \frac{1}{2}\mathbf{u}_i \mathbf{u}_i^\top + \frac{1}{2} \sum_{j=1}^i \sum_{k=0}^{j-1} h_{j,k}(\mathbf{u}_i \mathbf{u}_k^\top + \mathbf{u}_k \mathbf{u}_i^\top), \end{cases} \quad (22)$$

which results from the inequalities (19) at the points $(\mathbf{x}, \mathbf{y}) = (\mathbf{y}_{i-1}, \mathbf{y}_i)$ and $(\mathbf{x}, \mathbf{y}) = (\mathbf{x}_*, \mathbf{y}_i)$ respectively.

As in [14, Problem (DQ')], problem (P1) has a dual formulation that one could solve numerically using a semidefinite program (SDP) to determine an upper bound on the cost function convergence rate for any FSFOM:⁶

$$\begin{aligned} F(\mathbf{x}_N) - F(\mathbf{x}_*) &\leq \mathcal{B}_P(\mathbf{h}, N, d, L, R) \\ &\leq \mathcal{B}_D(\mathbf{h}, N, L, R) := \min_{\substack{(\boldsymbol{\lambda}, \boldsymbol{\tau}) \in A, \\ \boldsymbol{\gamma} \in \mathbb{R}}} \left\{ \frac{1}{2}LR^2\boldsymbol{\gamma} : \begin{pmatrix} \mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) & \frac{1}{2}\boldsymbol{\tau} \\ \frac{1}{2}\boldsymbol{\tau}^\top & \frac{1}{2}\boldsymbol{\gamma} \end{pmatrix} \succeq 0 \right\}, \end{aligned} \quad (D)$$

where $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_{N-1}]^\top \in \mathbb{R}_+^{N-1}$, $\boldsymbol{\tau} = [\tau_0, \dots, \tau_{N-1}]^\top \in \mathbb{R}_+^N$, and

$$A = \left\{ (\boldsymbol{\lambda}, \boldsymbol{\tau}) \in \mathbb{R}_+^{2N-1} : \begin{array}{l} \tau_0 = \lambda_1, \quad \lambda_{N-1} + \tau_{N-1} = 1, \\ \lambda_i - \lambda_{i+1} + \tau_i = 0, \quad i = 1, \dots, N-2, \end{array} \right\}, \quad (23)$$

$$\mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) = \sum_{i=1}^{N-1} \lambda_i \tilde{\mathbf{A}}_{i-1,i}(\mathbf{h}) + \sum_{i=0}^{N-1} \tau_i \tilde{\mathbf{D}}_i(\mathbf{h}). \quad (24)$$

This means that we could compute a valid upper bound (D) of (P) for given step coefficients \mathbf{h} using a SDP. The next two sections provide an analytical solution to (D) for FPGM and similarly for our new generalized FPGM, superseding the use of numerical SDP solvers.

3.2 Generalized FPGM

We specify a feasible point of (D) that leads to our new generalized form of FPGM.

Lemma 2 *For the following step coefficients:*

$$h_{i+1,k} = \begin{cases} \frac{t_{i+1}}{T_{i+1}} \left(t_k - \sum_{j=k+1}^i h_{j,k} \right), & k = 0, \dots, i-1, \\ 1 + \frac{(t_i-1)t_{i+1}}{T_{i+1}}, & k = i, \end{cases} \quad (25)$$

the choice of variables:

$$\lambda_i = \frac{T_{i-1}}{T_{N-1}}, \quad i = 1, \dots, N-1, \quad \tau_i = \frac{t_i}{T_{N-1}}, \quad i = 0, \dots, N-1, \quad \gamma = \frac{1}{T_{N-1}}, \quad (26)$$

is a feasible point of (D) for any choice of t_i such that

$$t_0 = 1, \quad t_i > 0, \quad \text{and} \quad t_i^2 \leq T_i := \sum_{l=0}^i t_l. \quad (27)$$

⁶ See Appendix for the derivation of the dual formulation (D) of (P1).

Proof It is obvious that $(\boldsymbol{\lambda}, \boldsymbol{\tau})$ in (26) with (27) is in Λ (23). Using (22), the (i, k) th entry of the symmetric matrix $\mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau})$ in (24) can be written as

$$S_{i,k}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) = \begin{cases} \frac{1}{2}((\lambda_i + \tau_i)h_{i,k} + \tau_i \sum_{j=k+1}^{i-1} h_{j,k}), & i = 2, \dots, N-1, k = 0, \dots, i-2, \\ \frac{1}{2}((\lambda_i + \tau_i)h_{i,i-1} - \lambda_i), & i = 1, \dots, N-1, k = i-1, \\ \frac{1}{2}\lambda_{i+1}, & i = 0, \dots, N-2, k = i, \\ \frac{1}{2}, & i = N-1, k = i, \end{cases}$$

where each element $S_{i,k}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau})$ corresponds to the coefficient of the term $\mathbf{u}_i \mathbf{u}_k^\top$ of $\mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau})$ in (24). Then, inserting (25) and (26) to the above yields

$$\begin{aligned} S_{i,k}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) &= \begin{cases} \frac{1}{2} \left(\frac{T_i}{T_{N-1}} \frac{t_i}{T_i} \left(t_k - \sum_{j=k+1}^{i-1} h_{j,k} \right) + \frac{t_i}{T_{N-1}} \sum_{j=k+1}^{i-1} h_{j,k} \right), & i = 2, \dots, N-1, k = 0, \dots, i-2, \\ \frac{1}{2} \left(\frac{T_i}{T_{N-1}} \left(1 + \frac{(t_{i-1}-1)t_i}{T_i} \right) - \frac{T_{i-1}}{T_{N-1}} \right), & i = 1, \dots, N-1, k = i-1, \\ \frac{T_i}{2T_{N-1}}, & i = 0, \dots, N-1, k = i. \end{cases} \\ &= \begin{cases} \frac{t_i t_k}{2T_{N-1}}, & i = 1, \dots, N-1, k = 0, \dots, i-1, \\ \frac{T_i}{2T_{N-1}}, & i = 0, \dots, N-1, k = i. \end{cases} \end{aligned}$$

Then, using (26) and (27), we finally show the feasibility condition of (D):

$$\begin{pmatrix} \mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) & \frac{1}{2}\boldsymbol{\tau} \\ \frac{1}{2}\boldsymbol{\tau}^\top & \frac{1}{2}\gamma \end{pmatrix} = \frac{1}{2T_{N-1}} (\text{diag}\{\mathbf{T} - \mathbf{t}^2\} + \mathbf{t}\mathbf{t}^\top) \succeq 0,$$

where $\mathbf{t} = (t_0, \dots, t_{N-1}, 1)^\top$ and $\mathbf{T} = (T_0, \dots, T_{N-1}, 1)^\top$.

FSFOM with the step coefficients (25) would be both computationally and memory-wise inefficient, so we next present an efficient equivalent recursive form of FSFOM with (25), named Generalized FPGM (GFPGM).

Algorithm GFPGM

Input: $f \in \mathcal{F}_L(\mathbb{R}^d)$, $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{y}_0 = \mathbf{x}_0$, $t_0 = T_0 = 1$.

For $i = 0, \dots, N-1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i)$$

Choose t_{i+1} s.t. $t_{i+1} > 0$ and $t_{i+1}^2 \leq T_{i+1} := \sum_{l=0}^{i+1} t_l$

$$\mathbf{y}_{i+1} = \mathbf{x}_{i+1} + \frac{(T_i - t_i)t_{i+1}}{t_i T_{i+1}}(\mathbf{x}_{i+1} - \mathbf{x}_i) + \frac{(t_i^2 - T_i)t_{i+1}}{t_i T_{i+1}}(\mathbf{x}_{i+1} - \mathbf{y}_i)$$

Proposition 1 The sequence $\{\mathbf{x}_0, \dots, \mathbf{x}_N\}$ generated by FSFOM with step sizes (25) is identical to the corresponding sequence generated by GFPGM.

Proof See Appendix.

Using Lemma 2, the following theorem bounds the cost function for the GFPGM iterates.

Theorem 3 Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $\mathcal{F}_L(\mathbb{R}^d)$ and let $\mathbf{x}_0, \dots, \mathbf{x}_N \in \mathbb{R}^d$ be generated by GFPGM. Then for $N \geq 1$,

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \frac{LR^2}{2T_{N-1}}. \quad (28)$$

Proof Using (D), Lemma 2 and Prop. 1, we have

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \mathcal{B}_D(\mathbf{h}, N, L, R) = \frac{1}{2}LR^2\gamma = \frac{LR^2}{2T_{N-1}}. \quad (29)$$

The GFPGM and Thm. 3 reduce to FPGM and (12) when $t_i^2 = T_i$ for all i , and Sec. 3.4 describes that FPGM results from optimizing the step coefficients of FSFOM with respect to the cost function form of the relaxed PEP (D). This GFPGM also includes the choice $t_i = \frac{i+a}{a}$ for any $a \geq 2$ as used in [8], which we denote as FPGM- a that differs from the algorithm in [8]. The following corollary provides a cost function convergence bound for FPGM- a .

Corollary 1 *Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $\mathcal{F}_L(\mathbb{R}^d)$ and let $\mathbf{x}_0, \dots, \mathbf{x}_N \in \mathbb{R}^d$ be generated by GFPGM with $t_i = \frac{i+a}{a}$ (FPGM- a) for any $a \geq 2$. Then for $N \geq 1$,*

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \frac{aLR^2}{N(N+2a-1)}. \quad (30)$$

Proof Thm. 3 implies (30), since $t_i = \frac{i+a}{a}$ satisfies (27), i.e.,

$$T_i - t_i^2 = \frac{(i+1)(i+2a)}{2a} - \frac{(i+a)^2}{a^2} = \frac{(a-2)i^2 + a(2a-3)i}{2a^2} \geq 0 \quad (31)$$

for any $a \geq 2$ and all $i \geq 0$.

3.3 Related work of GFPGM

This section shows that the GFPGM has a close connection to the accelerated algorithm in [27] that was developed specifically for a constrained smooth convex problem with a closed convex set Q , i.e.,

$$\phi(\mathbf{x}) = \mathbf{I}_Q(\mathbf{x}) := \begin{cases} 0, & \mathbf{x} \in Q, \\ \infty, & \text{otherwise.} \end{cases} \quad (32)$$

The projection operator $\mathbf{P}_Q(\mathbf{x}) := \arg \min_{\mathbf{y} \in Q} \|\mathbf{x} - \mathbf{y}\|$ is used for the proximal gradient update (2).

We show that the GFPGM can be written in the following equivalent form, named GFPGM', which is similar to that of the accelerated algorithm in [27] shown below. Note that the accelerated algorithm in [27] satisfies the bound (28) of the GFPGM in [27, Thm. 2] when $\phi(\mathbf{x}) = \mathbf{I}_Q(\mathbf{x})$.

Algorithm GFPGM'

Input: $f \in \mathcal{F}_L(\mathbb{R}^d)$, $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{y}_0 = \mathbf{x}_0$, $t_0 = T_0 = 1$.

For $i = 0, \dots, N-1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i) = \mathbf{y}_i - \frac{1}{L} \tilde{\nabla}_L F(\mathbf{y}_i)$$

$$\mathbf{z}_{i+1} = \mathbf{y}_0 - \frac{1}{L} \sum_{k=0}^i t_k \tilde{\nabla}_L F(\mathbf{y}_k)$$

$$\text{Choose } t_{i+1} \text{ s.t. } t_{i+1} > 0 \text{ and } t_{i+1}^2 \leq T_{i+1} = \sum_{l=0}^{i+1} t_l$$

$$\mathbf{y}_{i+1} = \left(1 - \frac{t_{i+1}}{T_{i+1}}\right) \mathbf{x}_{i+1} + \frac{t_{i+1}}{T_{i+1}} \mathbf{z}_{i+1}$$

Algorithm [27] for $\phi(\mathbf{x}) = \mathbf{I}_Q(\mathbf{x})$

Input: $f \in \mathcal{F}_L(\mathbb{R}^d)$, $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{y}_0 = \mathbf{x}_0$, $t_0 = T_0 = 1$.

For $i = 0, \dots, N-1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i) = \mathbf{P}_Q \left(\mathbf{y}_i - \frac{1}{L} \nabla f(\mathbf{y}_i) \right)$$

$$\mathbf{z}_{i+1} = \mathbf{P}_Q \left(\mathbf{y}_0 - \frac{1}{L} \sum_{k=0}^i t_k \nabla f(\mathbf{y}_k) \right)$$

$$\text{Choose } t_{i+1} \text{ s.t. } t_{i+1} > 0 \text{ and } t_{i+1}^2 \leq T_{i+1} = \sum_{l=0}^{i+1} t_l$$

$$\mathbf{y}_{i+1} = \left(1 - \frac{t_{i+1}}{T_{i+1}} \right) \mathbf{x}_{i+1} + \frac{t_{i+1}}{T_{i+1}} \mathbf{z}_{i+1}$$

Proposition 2 *The sequence $\{\mathbf{x}_0, \dots, \mathbf{x}_N\}$ generated by GFPGM is identical to the corresponding sequence generated by GFPGM'.*

Proof See Appendix.

Clearly GFPGM' and the accelerated algorithm in [27] are equivalent for the unconstrained smooth convex problem ($Q = \mathbb{R}^d$). However, when the operation $\mathbf{P}_Q(\mathbf{x})$ is relatively expensive, our GFPGM and GFPGM' that use one projection per iteration could be preferred over the accelerated algorithm in [27] that uses two projections per iteration.

3.4 Optimizing step coefficients of FSFOM using the cost function form of PEP

To find the step coefficients in the class FSFOM that are optimal in terms of the cost function form of PEP, we would like to solve the following problem:

$$\hat{\mathbf{h}}_{\text{P}} := \arg \min_{\mathbf{h} \in \mathbb{R}^{N(N+1)/2}} \mathcal{B}_{\text{P}}(\mathbf{h}, N, d, L, R). \quad (\text{HP})$$

Because (HP) seems intractable, we instead optimize the step coefficients using the relaxed bound in (D):

$$\hat{\mathbf{h}}_{\text{D}} := \arg \min_{\mathbf{h} \in \mathbb{R}^{N(N+1)/2}} \mathcal{B}_{\text{D}}(\mathbf{h}, N, L, R). \quad (\text{HD})$$

The problem (HD) is bilinear, and a convex relaxation technique in [14, Thm. 3] makes it efficiently solvable using numerical methods. We optimized (HD) numerically for many choices of N using a SDP solver [10, 18] and based on our numerical results (not shown) we conjecture that the feasible point in Lemma 2 with $t_i^2 = T_i$ that corresponds to FPGM (FISTA) is a global minimizer of (HD). It is straightforward to show that the step coefficients in Lemma 2 with $t_i^2 = T_i$ give the smallest bound of (D) and (28) among all feasible points in Lemma 2, but showing optimality among all possible feasible points of (HD) may require further derivations as in [20, Lemma 3] using KKT conditions, which we leave as future work.

This section has provided an alternate convergence proof of FPGM using the new relaxed PEP, and suggested that FPGM corresponds to FSFOM with optimized step coefficients using the cost function form of the relaxed PEP. Because minimizing the norm of the (composite) gradient (mapping) could be important in addition to the cost function in dual problems, the next section provides an alternate optimization of the step coefficients of FSFOM with respect to the norm of the composite gradient mapping.

4 Relaxation and optimization of the composite gradient mapping form of PEP

4.1 Relaxation for the composite gradient mapping form of PEP

To form a worst-case convergence bound on the norm of the composite gradient mapping for a given \mathbf{h} of FSFOM, we use the following PEP that replaces $F(\mathbf{x}_N) - F(\mathbf{x}_*)$ in (P) by the norm squared of the composite gradient mapping. Here, we consider the smallest composite gradient mapping norm squared among all iterates⁷ ($\min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|L(\mathbf{p}_L(\mathbf{x}) - \mathbf{x})\|^2 = \min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\|^2$) as follows:

$$\begin{aligned} \mathcal{B}_{P'}(\mathbf{h}, N, d, L, R) := & \max_{\substack{F \in \mathcal{F}_L(\mathbb{R}^d), \\ \mathbf{x}_0, \dots, \mathbf{x}_N, \mathbf{x}_* \in \mathbb{R}^d, \\ \mathbf{y}_0, \dots, \mathbf{y}_{N-1} \in \mathbb{R}^d, \\ \mathbf{x}_* \in X_*(F), \|\mathbf{x}_0 - \mathbf{x}_*\| \leq R}} \min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|L(\mathbf{p}_L(\mathbf{x}) - \mathbf{x})\|^2 \\ \text{s.t. } & \mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i), \quad i = 0, \dots, N-1, \\ & \mathbf{y}_{i+1} = \mathbf{y}_i + \sum_{k=0}^i h_{i+1,k}(\mathbf{x}_{k+1} - \mathbf{y}_k), \quad i = 0, \dots, N-2, \end{aligned} \quad (P')$$

Because this infinite-dimensional max-min problem appears intractable, similar to the relaxation from (P) to (P1), we relax (P') to a finite-dimensional problem using an additional constraint resulting from (6) that is equivalent to

$$\frac{L}{2} \|\mathbf{p}_L(\mathbf{x}_N) - \mathbf{x}_N\|^2 \leq F(\mathbf{x}_N) - F(\mathbf{x}_*) \quad (33)$$

and conditions that are equivalent to $\alpha \leq \|L(\mathbf{p}_L(\mathbf{x}) - \mathbf{x})\|^2$ for all $\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}$ after replacing $\min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|L(\mathbf{p}_L(\mathbf{x}) - \mathbf{x})\|^2$ by α as in [33]:

$$\begin{aligned} \mathcal{B}_{P1'}(\mathbf{h}, N, d, L, R) := & \max_{\substack{\bar{\mathbf{G}} \in \mathbb{R}^{(N+1) \times d}, \\ \bar{\boldsymbol{\delta}} \in \mathbb{R}^N, \alpha \in \mathbb{R}}} L^2 R^2 \alpha \\ \text{s.t. } & \text{Tr}\{\bar{\mathbf{G}}^\top \bar{\mathbf{A}}_{i-1,i}(\mathbf{h}) \bar{\mathbf{G}}\} \leq \delta_{i-1} - \delta_i, \quad i = 1, \dots, N-1, \\ & \text{Tr}\{\bar{\mathbf{G}}^\top \bar{\mathbf{D}}_i(\mathbf{h}) \bar{\mathbf{G}} + \nu \bar{\mathbf{u}}_i^\top \bar{\mathbf{G}}\} \leq -\delta_i, \quad i = 0, \dots, N-1, \\ & \text{Tr}\left\{\frac{1}{2} \bar{\mathbf{G}}^\top \bar{\mathbf{u}}_N \bar{\mathbf{u}}_N^\top \bar{\mathbf{G}}\right\} \leq \delta_{N-1}, \\ & \text{Tr}\{-\bar{\mathbf{G}}^\top \bar{\mathbf{u}}_i \bar{\mathbf{u}}_i^\top \bar{\mathbf{G}}\} \leq -\alpha, \quad i = 0, \dots, N, \end{aligned} \quad (P1')$$

for any given unit vector $\boldsymbol{\nu} \in \mathbb{R}^d$, by defining the $(i+1)$ th standard basis vector $\bar{\mathbf{u}}_i = \mathbf{e}_{i+1} \in \mathbb{R}^{N+1}$, the matrices

$$\bar{\mathbf{A}}_{i-1,i}(\mathbf{h}) = \begin{pmatrix} \check{\mathbf{A}}_{i-1,i}(\mathbf{h}) & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix}, \quad \bar{\mathbf{D}}_i(\mathbf{h}) = \begin{pmatrix} \check{\mathbf{D}}_i(\mathbf{h}) & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix} \quad (34)$$

where $\mathbf{0} = [0, \dots, 0]^\top \in \mathbb{R}^N$, and the matrix $\bar{\mathbf{G}} = [\mathbf{G}^\top, \bar{\mathbf{g}}_N]^\top \in \mathbb{R}^{(N+1) \times d}$ where

$$\bar{\mathbf{g}}_N := -\frac{1}{\|\mathbf{y}_0 - \mathbf{x}_*\|} (\mathbf{p}_L(\mathbf{x}_N) - \mathbf{x}_N) = \frac{1}{L\|\mathbf{y}_0 - \mathbf{x}_*\|} \tilde{\nabla}_L F(\mathbf{x}_N). \quad (35)$$

Similar to (D) and [19, Problem (D'')], we have the following dual formulation of (P1') that could be solved using SDP:

$$\mathcal{B}_{D'}(\mathbf{h}, N, L, R) := \min_{\substack{(\boldsymbol{\lambda}, \boldsymbol{\tau}, \boldsymbol{\eta}, \boldsymbol{\beta}) \in \mathcal{A}', \\ \gamma \in \mathbb{R}}} \left\{ \frac{1}{2} L^2 R^2 \gamma : \begin{pmatrix} \mathbf{S}'(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}, \boldsymbol{\eta}, \boldsymbol{\beta}) & \frac{1}{2} [\boldsymbol{\tau}^\top, 0]^\top \\ \frac{1}{2} [\boldsymbol{\tau}^\top, 0] & \frac{1}{2} \gamma \end{pmatrix} \succeq 0 \right\} \quad (D')$$

⁷ See Appendix.

where $\eta \in \mathbb{R}_+$, $\beta = [\beta_0, \dots, \beta_N]^\top \in \mathbb{R}_+^{N+1}$, and

$$\Lambda' = \left\{ (\lambda, \tau, \eta, \beta) \in \mathbb{R}_+^{3N+1} : \begin{array}{l} \tau_0 = \lambda_1, \quad \lambda_{N-1} + \tau_{N-1} = \eta, \quad \sum_{i=0}^N \beta_i = 1, \\ \lambda_i - \lambda_{i+1} + \tau_i = 0, \quad i = 1, \dots, N-2 \end{array} \right\}, \quad (36)$$

$$\mathbf{S}'(\mathbf{h}, \lambda, \tau, \eta, \beta) = \sum_{i=1}^{N-1} \lambda_i \bar{\mathbf{A}}_{i-1,i}(\mathbf{h}) + \sum_{i=0}^{N-1} \tau_i \bar{\mathbf{D}}_i(\mathbf{h}) + \frac{1}{2} \eta \bar{\mathbf{u}}_N \bar{\mathbf{u}}_N^\top - \sum_{i=0}^N \beta_i \bar{\mathbf{u}}_i \bar{\mathbf{u}}_i^\top. \quad (37)$$

The next section specifies a feasible point of interest that is in the class of GFPGM and analyzes the convergence rate of the norm of the composite gradient mapping. Then we optimize the step coefficients of FSFOM with respect to the composite gradient mapping form of PEP leading to a new algorithm that differs from Nesterov's acceleration for decreasing the cost function.

4.2 Convergence analysis of the composite gradient mapping of GFPGM

The following Lemma provides feasible point of (D') for the step coefficients (25) of GFPGM.

Lemma 3 *For the step coefficients $\{h_{i+1,k}\}$ in (25), the choice of variables*

$$\lambda_i = T_{i-1} \tau_0, \quad i = 1, \dots, N-1, \quad \tau_i = \begin{cases} \left(\frac{1}{2} \left(\sum_{k=0}^{N-1} (T_k - t_k^2) + T_{N-1} \right) \right)^{-1}, & i = 0, \\ t_i \tau_0, & i = 1, \dots, N-1, \end{cases} \quad (38)$$

$$\eta = T_{N-1} \tau_0, \quad \beta_i = \begin{cases} \frac{1}{2} (T_i - t_i^2) \tau_0, & i = 0, \dots, N-1, \\ \frac{1}{2} T_{N-1} \tau_0, & i = N, \end{cases} \quad \gamma = \tau_0. \quad (39)$$

is a feasible point of (D') for any choice of t_i and T_i satisfying (27).

Proof It is obvious that $(\lambda, \tau, \eta, \beta)$ in (38) and (39) with (27) is in Λ' (36). Using (22) and (34), the (i, k) th entry of the symmetric matrix $\mathbf{S}'(\mathbf{h}, \lambda, \tau, \eta, \beta)$ in (37) can be written as

$$\begin{aligned} & S'_{i,k}(\mathbf{h}, \lambda, \tau, \eta, \beta) \\ &= \begin{cases} \frac{1}{2} \left((\lambda_i + \tau_i) h_{i,k} + \tau_i \sum_{j=k+1}^{i-1} h_{j,k} \right), & i = 2, \dots, N-1, \quad k = 0, \dots, i-2, \\ \frac{1}{2} ((\lambda_i + \tau_i) h_{i,i-1} - \lambda_i), & i = 1, \dots, N-1, \quad k = i-1, \\ \frac{1}{2} \lambda_{i+1} - \beta_i, & i = 0, \dots, N-2, \quad k = i, \\ \frac{1}{2} \eta - \beta_i, & i = N-1, N, \quad k = i, \\ 0, & i = N, \quad k = 0, \dots, i-1, \end{cases} \end{aligned}$$

and inserting (25), (38), and (39) yields

$$\begin{aligned} & S'_{i,k}(\mathbf{h}, \lambda, \tau, \eta, \beta) \\ &= \begin{cases} \frac{1}{2} \left(T_i \tau_0 \frac{t_i}{T_i} \left(t_k - \sum_{j=k+1}^{i-1} h_{j,k} \right) + t_i \tau_0 \sum_{j=k+1}^{i-1} h_{j,k} \right), & i = 2, \dots, N-1, \quad k = 0, \dots, i-2, \\ \frac{1}{2} \left(T_i \tau_0 \left(1 + \frac{(t_{i-1}-1)t_i}{T_i} \right) - T_{i-1} \tau_0 \right), & i = 1, \dots, N-1, \quad k = i-1, \\ \frac{1}{2} T_i \tau_0 - \frac{1}{2} (T_i - t_i^2) \tau_0, & i = 0, \dots, N-1, \quad k = i, \\ 0, & i = N, \quad k = 0, \dots, i, \end{cases} \\ &= \begin{cases} \frac{1}{2} t_i t_k \tau_0, & i = 0, \dots, N-1, \quad k = 0, \dots, i, \\ 0, & i = N, \quad k = 0, \dots, i. \end{cases} \end{aligned}$$

Finally, by defining $\bar{\mathbf{t}} = (t_0, \dots, t_{N-1}, 0, 1)^\top$ we have the feasibility condition of (D'):

$$\begin{pmatrix} \mathbf{S}'(\mathbf{h}, \lambda, \tau, \eta, \beta) \frac{1}{2} [\bar{\mathbf{t}}^\top, 0]^\top \\ \frac{1}{2} [\bar{\mathbf{t}}^\top, 0] \quad \frac{1}{2} \gamma \end{pmatrix} = \frac{1}{2} \bar{\mathbf{t}} \bar{\mathbf{t}}^\top \tau_0 \succeq 0.$$

Using Lemma 3, the following theorem bounds the (smallest) norm of the composite gradient mapping for the GFPGM iterates.

Theorem 4 *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $\mathcal{F}_L(\mathbb{R}^d)$ and let $\mathbf{x}_0, \dots, \mathbf{x}_N, \mathbf{y}_0, \dots, \mathbf{y}_{N-1} \in \mathbb{R}^d$ be generated by GFPGM. Then for $N \geq 1$,*

$$\min_{\mathbf{x} \in \{\mathbf{x}_0, \dots, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\| \leq \min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\| \leq \frac{LR}{\sqrt{\sum_{k=0}^{N-1} (T_k - t_k^2) + T_{N-1}}}. \quad (40)$$

Proof Lemma 1 implies the first inequality of (40). Using (D'), Lemma 3 and Prop. 1, we have

$$\begin{aligned} \min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\|^2 &\leq \mathcal{B}_{D'}(\mathbf{h}, N, L, R) = \frac{1}{2} L^2 R^2 \gamma \\ &= \frac{L^2 R^2}{\sum_{k=0}^{N-1} (T_k - t_k^2) + T_{N-1}}, \end{aligned}$$

which is equivalent to (40).

Although the bound (40) is not guaranteed to be tight due to the relaxation on PEP, next two sections show that there exists choices of t_i that provide a rate $O(1/N^{\frac{3}{2}})$ for decreasing the composite gradient mapping, including the choice that optimizes the composite gradient mapping form of PEP.

FGM for smooth convex minimization was shown to achieve the rate $O(1/N^{\frac{3}{2}})$ for the decrease of the usual gradient in [19]. In contrast, Thm. 4 provides only a $O(1/N)$ bound for FPGM (or GFPGM with t_i (11)) on the decrease of the composite gradient mapping since $T_i = t_i^2$ for all i and the value of T_{N-1} is $O(N^2)$ for t_i (11). Sec. 5 below numerically studies a tight bound on the composite gradient mapping of FPGM and illustrates that it has a rate that is faster than the rate $O(1/N)$ of Thm. 4, indicating there is a room for improvement in the composite gradient mapping form of the relaxed PEP

4.3 Optimizing step coefficients of FSFOM using the composite gradient mapping form of PEP

To optimize the step coefficients in the class FSFOM in terms of the composite gradient mapping form of the relaxed PEP (D'), we would like to solve the following problem:

$$\hat{\mathbf{h}}_{D'} := \arg \min_{\mathbf{h} \in \mathbb{R}^{N(N+1)/2}} \mathcal{B}_{D'}(\mathbf{h}, N, L, R). \quad (\text{HD}')$$

Similar to (HD), we use a convex relaxation [14, Thm. 3] to make the bilinear problem (HD') efficiently solvable using numerical methods. We then numerically optimized (HD') for many choices of N using a SDP solver [10, 18] and found that the following choice of t_i :

$$t_i = \begin{cases} 1, & i = 0, \\ \frac{1 + \sqrt{1 + 4t_{i-1}^2}}{2}, & i = 1, \dots, \lfloor \frac{N}{2} \rfloor - 1, \\ \frac{N-i+1}{2}, & i = \lfloor \frac{N}{2} \rfloor, \dots, N-1, \end{cases} \quad (41)$$

makes the feasible point in Lemma 3 optimal empirically with respect to the relaxed bound (HD'). Interestingly, whereas the usual t_i factors increase with i indefinitely, here, the factors begin decreasing after $i = \lfloor \frac{N}{2} \rfloor - 1$.

We also noticed numerically that finding the t_i that minimizes the bound (40), *i.e.*, solving the following constrained quadratic problem:

$$\max_{\{t_i\}} \left\{ \sum_{k=0}^{N-1} \left(\sum_{l=0}^k t_l - t_k^2 \right) + \sum_{l=0}^{N-1} t_l \right\} \quad \text{s.t.} \quad t_i \text{ satisfies (27) for all } i, \quad (42)$$

is equivalent to optimizing (HD'). This means that the solution of (42) numerically appears equivalent to (41), the (conjectured) solution of (HD'). Interestingly, the unconstrained maximizer of (42) without the constraint (27) is $t_i = \frac{N-i+1}{2}$, and this partially appears in the constrained maximizer (41) of the problem (42).

Based on this numerical evidence, we conjecture that the solution $\hat{\mathbf{h}}_{D'}$ of problem (HD') corresponds to (25) with (41). Using Prop. 1, the following GFPGM form with (41) is equivalent to FSFOM with the step coefficients (25) for (41) that are optimized step coefficients of FSFOM with respect to the norm of the composite gradient mapping, which we name FPGM-OCG (OCG for optimized over composite gradient mapping).

Algorithm FPGM-OCG (GFPGM with t_i in (41))

Input: $f \in C_L^{1,1}(\mathbb{R}^d)$ convex, $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{y}_0 = \mathbf{x}_0$, $t_0 = T_0 = 1$.

For $i = 0, \dots, N-1$

$$\mathbf{x}_{i+1} = \mathbf{p}_L(\mathbf{y}_i)$$

$$t_{i+1} = \begin{cases} \frac{1+\sqrt{1+4t_i^2}}{2}, & i = 1, \dots, \lfloor \frac{N}{2} \rfloor - 2, \\ \frac{N-i}{2}, & i = \lfloor \frac{N}{2} \rfloor - 1, \dots, N-2, \end{cases}$$

$$\begin{aligned} \mathbf{y}_{i+1} = \mathbf{x}_{i+1} &+ \frac{(T_i - t_i)t_{i+1}}{t_i T_{i+1}}(\mathbf{x}_{i+1} - \mathbf{x}_i) \\ &+ \frac{(t_i^2 - T_i)t_{i+1}}{t_i T_{i+1}}(\mathbf{x}_{i+1} - \mathbf{y}_i), \quad i < N-1 \end{aligned}$$

The following theorem bounds the cost function and the (smallest) norm of the composite gradient mapping for the FPGM-OCG iterates.

Theorem 5 Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $\mathcal{F}_L(\mathbb{R}^d)$ and let $\mathbf{x}_0, \dots, \mathbf{x}_N, \mathbf{y}_0, \dots, \mathbf{y}_{N-1} \in \mathbb{R}^d$ be generated by FPGM-OCG. Then for $N \geq 1$,

$$F(\mathbf{x}_N) - F(\mathbf{x}_*) \leq \frac{4L\|\mathbf{x}_0 - \mathbf{x}_*\|^2}{N(N+4)}, \quad (43)$$

and for $N \geq 3$,

$$\min_{\mathbf{x} \in \{\mathbf{x}_0, \dots, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\| \leq \min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\| \leq \frac{2\sqrt{6}LR}{N\sqrt{N-2}}. \quad (44)$$

Proof FPGM-OCG is an instance of the GFPGM, and thus Thm. 3 implies (43) using

$$\begin{aligned} T_{N-1} &= T_{m-1} + \sum_{k=m}^{N-1} t_k = t_{m-1}^2 + \sum_{k=m}^{N-1} \frac{N-k+1}{2} = t_{m-1}^2 + \sum_{k'=2}^{N-m+1} \frac{k'}{2} \\ &\geq \frac{(m+1)^2}{4} + \frac{(N-m+1)(N-m+2)}{4} - \frac{1}{2} \geq \frac{2N^2 + 8N + 1}{16}, \end{aligned}$$

where $m = \lfloor \frac{N}{2} \rfloor \geq \frac{N-1}{2}$, $N-m \geq \frac{N}{2}$, and $T_{m-1} = t_{m-1}^2 \geq \frac{(m+1)^2}{4}$ (13).

In addition, Thm. 4 implies (44), using

$$\sum_{k=0}^{N-1} (T_k - t_k^2) + T_{N-1} \geq \frac{1}{24}(N-2)N^2, \quad (45)$$

which we prove in the Appendix.

The composite gradient mapping bound (44) of FPGM-OCG is asymptotically $\frac{2\sqrt{2}}{3}$ -times smaller than the bound (15) of FPGM- $(m = \lfloor \frac{2N}{3} \rfloor)$. In addition, the cost function bound (43) of FPGM-OCG satisfies the optimal rate $O(1/N^2)$, while the bound (43) is two-times larger than the analogous bound (12) of FPGM.

4.4 Decreasing the composite gradient mapping with a rate $O(1/N^{\frac{3}{2}})$ without selecting N in advance

FPGM-OCG and FPGM- m satisfy a fast rate $O(1/N^{\frac{3}{2}})$ for decreasing the norm of the composite gradient mapping but require one to select the total number of iterations N in advance, which could be undesirable in practice. One could use FPGM- σ in [23] that does not require selecting N in advance, but instead we suggest a new choice of t_i in GFPGM that satisfies a composite gradient mapping bound that is lower than the bound (17) of FPGM- σ .

Based on Thm. 4, the following corollary shows that GFPGM with $t_i = \frac{i+a}{a}$ (FPGM- a) for any $a > 2$ satisfies the rate $O(1/N^{\frac{3}{2}})$ of the norm of the composite gradient mapping without selecting N in advance. (Cor. 1 showed that FPGM- a for any $a \geq 2$ satisfies the optimal rate $O(1/N^2)$ of the cost function.)

Corollary 2 *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be in $\mathcal{F}_L(\mathbb{R}^d)$ and let $\mathbf{x}_0, \dots, \mathbf{x}_N, \mathbf{y}_0, \dots, \mathbf{y}_{N-1} \in \mathbb{R}^d$ be generated by GFPGM with $t_i = \frac{i+a}{a}$ (FPGM- a) for any $a \geq 2$. Then for $N \geq 1$, we have the following bound on the (smallest) composite gradient mapping:*

$$\begin{aligned} \min_{\mathbf{x} \in \{\mathbf{x}_0, \dots, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\| &\leq \min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\| \\ &\leq \frac{a\sqrt{6}LR}{\sqrt{N((a-2)N^2 + 3(a^2 - a + 1)N + (3a^2 + 2a - 1))}}. \end{aligned} \quad (46)$$

Proof With $T_i = \frac{(i+1)(i+2a)}{2a}$ and (31), Thm. 4 implies (46) using

$$\begin{aligned} \sum_{k=0}^{N-1} (T_k - t_k^2) + T_{N-1} &= \sum_{k=0}^{N-1} \left(\frac{(k+1)(k+2a)}{2a} - \frac{(k+a)^2}{a^2} \right) + \frac{N(N+2a-1)}{2a} \\ &= \sum_{k=0}^{N-1} \left(\frac{(a-2)k^2 + a(2a-3)k}{2a^2} \right) + \frac{N(N+2a-1)}{2a} \\ &= \frac{N}{2a^2} \left(\frac{(a-2)(N-1)(2N-1)}{6} + \frac{a(2a-3)(N-1)}{2} + a(N+2a-1) \right) \\ &= \frac{N((a-2)N^2 + 3(a^2 - a + 1)N + (3a^2 + 2a - 1))}{6a^2}. \end{aligned}$$

FPGM- a for any $a > 2$ has a composite gradient mapping bound (46) that is asymptotically $\frac{a}{2\sqrt{a-2}}$ -times larger than the bound (44) of FPGM-OCG. This gap reduces to $\sqrt{2}$ at best when $a = 4$, which is clearly better than that of FPGM- σ . Therefore, this FPGM- a algorithm will be useful for minimizing the composite gradient mapping with a rate $O(1/N^{\frac{3}{2}})$ without selecting N in advance.

5 Discussion and Conclusion

Table 1 summarizes the asymptotic convergence rate bounds of all algorithms discussed in this paper. (Note that the bounds are not guaranteed to be tight.) In Table 1, FPGM and FPGM-OCG provide the best known analytical worst-case convergence bounds for decreasing the cost function and the composite gradient mapping respectively. When one does not want to select N in advance for decreasing the composite gradient mapping, FPGM- a will be a useful alternative to FPGM-OCG.

Since none of the bounds presented in Table 1 are guaranteed to be tight, we modified the code⁸ in Taylor *et al.* [32] to compare tight (numerical) bounds for the cost function and the composite gradient mapping in Tables 2 and 3 respectively for $N = 1, 2, 4, 10, 20, 30, 40, 47, 50$. This numerical bound is guaranteed to be tight when the large-scale condition is satisfied [32]. Taylor *et al.* [32, Fig. 1] already studied a tight worst-case bound on the cost function decrease of FPGM numerically, and found that the analytical bound (12) is asymptotically tight. Table 2 additionally provides numerical tight bounds on the cost function of all algorithms presented in this paper, also suggesting that our relaxation of the cost

⁸ The code in Taylor *et al.* [32] currently does not provide a tight bound of the norm of the composite gradient mapping, so we simply added few lines to compute a tight bound.

Algorithm	Asymptotic convergence rate bound		Require selecting N in advance
	Cost function ($\times LR^2$)	Proximal gradient ($\times LR$)	
PGM	$\frac{1}{2}N^{-1}$	$2N^{-1}$	No
FPGM	$2N^{-2}$	$2N^{-1}$	No
FPGM- σ ($0 < \sigma < 1$)	$\frac{2}{\sigma^2}N^{-2}$	$\frac{2\sqrt{3}}{\sigma^2}\sqrt{\frac{1+\sigma}{1-\sigma}}N^{-\frac{3}{2}}$	No
FPGM- $(\sigma=0.78)$	$3.3N^{-2}$	$16.2N^{-\frac{3}{2}}$	No
FPGM- $(m = \lfloor \frac{2N}{3} \rfloor)$	$4.5N^{-2}$	$5.2N^{-\frac{3}{2}}$	Yes
FPGM-OCG	$4N^{-2}$	$4.9N^{-\frac{3}{2}}$	Yes
FPGM-a ($a > 2$)	aN^{-2}	$\frac{a\sqrt{6}}{\sqrt{a-2}}N^{-\frac{3}{2}}$	No
FPGM-$(a=4)$	$4N^{-2}$	$6.9N^{-\frac{3}{2}}$	No

Table 1 Asymptotic convergence rate bounds on the cost function $F(\mathbf{x}_N) - F(\mathbf{x}_*)$ and the norm of the composite gradient mapping $\min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\|$ of PGM, FPGM, FPGM- σ , FPGM- m , FPGM-OCG, and FPGM- a . (The cost function bound for FPGM- m in the table corresponds to the bound for FPGM after m iterations because a tight bound for the final N th iteration is unknown. The bound on $\min_{i \in \{0, \dots, N\}} \|\tilde{\nabla}_{L/\sigma^2} F(\mathbf{y}_i)\|$ is used for FPGM- σ .)

function form of the PEP from (P) to (D) is asymptotically tight (for some algorithms). In addition, the trend of the tight bounds of the composite gradient mapping in Table 3 follows that of the bounds in Table 1. However, there is gap between them that is not asymptotically tight, unlike the gap between the bounds of the cost function in Tables 1 and 2. In particular, the numerical tight bound for the composite gradient mapping of FPGM in Table 3 has a rate faster than the known rate $O(1/N)$ in Thm. 4. We leave reducing this gap for the bounds on the norm of the composite gradient mapping as future work, possibly with a tighter relaxation of PEP. In addition, FPGM- $(m = \lfloor \frac{2N}{3} \rfloor)$ has a numerical tight bound in Table 3 that is even slightly better than that of FPGM-OCG, unlike our expectation from the analytical bounds in Table 1 and Sec. 4.3. This shows room for improvement in optimizing the step coefficients of FSFOM with respect to the composite gradient mapping, again possibly with a tighter relaxation of PEP.

N	PGM	FPGM	FPGM - $(\sigma=0.78)$	FPGM - $(m = \lfloor \frac{2N}{3} \rfloor)$	FPGM -OCG	FPGM - $(a=4)$
1	4.00	4.00	2.43	4.00	4.00	4.00
2	8.00	8.00	4.87	8.00	8.00	8.00
4	16.00	19.35	11.77	17.13	17.60	17.23
10	40.00	79.07	48.11	56.47	59.25	55.88
20	80.00	261.66	159.19	163.75	170.10	159.17
30	120.00	546.51	332.49	321.56	331.97	312.03
40	160.00	932.89	567.57	502.37	544.55	514.73
47	188.00	1263.58	768.76	675.68	723.06	686.33
50	200.00	1420.45	864.20	752.90	807.66	767.37
Empi. $O(\cdot)$	$N^{-1.00}$	$N^{-1.89}$	$N^{-1.89}$	$N^{-1.75}$	$N^{-1.79}$	$N^{-1.80}$
Known $O(\cdot)$	N^{-1}	N^{-2}	N^{-2}	N^{-2}	N^{-2}	N^{-2}

Table 2 Tight worst-case convergence bounds on the cost function $LR^2/(F(\mathbf{x}_N) - F(\mathbf{x}_*))$ of PGM, FPGM, FPGM- $(\sigma=0.78)$, FPGM- $(m = \lfloor \frac{2N}{3} \rfloor)$, FPGM-OCG, and FPGM- $(a=4)$. We computed empirical rates by assuming that the bounds follow the form bN^c with constants b and c , and then by estimating c from points $N = 47, 50$. Note that the corresponding empirical rates are underestimated due to the simplified exponential model.

This paper focused on analyzing the convergence rate of the *smallest* composite gradient mapping among all iterates ($\min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\|$) in addition to the cost function, whereas the composite gradient mapping at the *final* iterate ($\|\tilde{\nabla}_L F(\mathbf{x}_N)\|$) could be also considered (see Appendix). For example, the composite gradient mapping bounds (10) and (15) for PGM and FPGM- m also apply to the final composite gradient mapping, and using (6) we can easily derive a (loose) convergence bound on the final composite gradient mapping for other algorithms, *e.g.*, such a final composite gradient mapping

N	PGM	FPGM	FPGM -($\sigma=0.78$)	FPGM -($m=\lfloor \frac{2N}{3} \rfloor$)	FPGM -OCG	FPGM -($a=4$)
1	1.84	1.84	1.18	1.84	1.84	1.84
2	2.83	2.83	1.78	2.83	2.83	2.83
4	4.81	5.65	3.50	5.09	5.21	5.12
10	10.80	13.24	8.74	14.91	15.60	14.76
20	20.78	27.19	18.83	39.70	39.61	29.21
30	30.78	43.49	30.82	64.45	64.40	47.14
40	40.78	61.76	44.39	92.82	91.99	67.82
47	47.77	75.60	54.73	113.92	113.41	83.67
50	50.77	81.78	59.35	123.54	123.17	90.78
Empi. $O(\cdot)$	$N^{-0.98}$	$N^{-1.27}$	$N^{-1.31}$	$N^{-1.31}$	$N^{-1.33}$	$N^{-1.32}$
Known $O(\cdot)$	N^{-1}	N^{-1}	$N^{-\frac{3}{2}}$	$N^{-\frac{3}{2}}$	$N^{-\frac{3}{2}}$	$N^{-\frac{3}{2}}$

Table 3 Tight worst-case convergence bounds on the norm of the composite gradient mapping $LR/(\min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L F(\mathbf{x})\|)$ of PGM, FPGM, FPGM-($\sigma=0.78$), FPGM-($m=\lfloor \frac{2N}{3} \rfloor$), FPGM-OCG, and FPGM-($a=4$). Empirical rates were computed as described in Table 2. (The bound for FPGM- σ uses $\min_{\mathbf{x} \in \{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}} \|\tilde{\nabla}_L/\sigma^2 F(\mathbf{x})\|$.)

bound for GFPGM is as follows:

$$\|\tilde{\nabla}_L F(\mathbf{x}_N)\| \stackrel{(6)}{\leq} \sqrt{2L(F(\mathbf{x}_N) - F(\mathbf{p}_L(\mathbf{x}_N)))} \leq \sqrt{2L(F(\mathbf{x}_N) - F(\mathbf{x}_*))} \stackrel{(28)}{\leq} \frac{LR}{\sqrt{T_{N-1}}}. \quad (47)$$

Since the optimal rate for decreasing the cost function is $O(1/N^2)$, the composite gradient mapping convergence bound (e.g., (47)) that is derived using (6) can provide only a rate $O(1/N)$ at best. To best of our knowledge, FPGM- m (or algorithms that similarly perform accelerated algorithms in the beginning and run a gradient method for the remaining iterations) is known only to have a rate $O(1/N^{\frac{3}{2}})$ in (15) for decreasing the final composite gradient mapping. Therefore, searching for first-order methods that have a convergence bound on the final composite gradient mapping that is lower than that of FPGM- m and possibly do not require knowing N in advance is an interesting open problem. Note that a regularization technique in [28] that provides a faster rate $O(1/N^2)$ (up to a logarithmic factor) for decreasing the final gradient norm for smooth convex minimization can be easily extended for rapidly minimizing the final composite gradient mapping with such rate for the composite problem (M); however, that approach requires knowing R in advance.

In conclusion, this paper analyzed and developed fixed-step first-order methods (FSFOM) for nonsmooth composite convex cost functions. We showed an alternate proof of FPGM (FISTA) using PEP, and suggested that FPGM (FISTA) results from optimizing the step coefficients of FSFOM with respect to the cost function form of the (relaxed) PEP. We then described a new generalized version of FPGM and analyzed its convergence using the (relaxed) PEP over both the cost function and the norm of the composite gradient mapping. Furthermore, we optimized the step coefficients of FSFOM with respect to the composite gradient mapping form of the (relaxed) PEP, yielding FPGM-OCG, which could be useful particularly when tackling dual problems.

Our relaxed PEP provided tractable and interesting analysis of the optimized step coefficients of FSFOM with respect to the cost function and the norm of the composite gradient mapping, but the relaxation is not guaranteed to be tight and the corresponding accelerations of PGM (FPGM and FPGM-OCG) are thus unlikely to be optimal. Therefore, finding an optimal step coefficients of FSFOM over the cost function and the norm of the composite gradient mapping remain as future work. Nevertheless, the proposed FPGM-OCG that optimizes the composite gradient mapping form of the relaxed PEP and the FPGM- a (for any $a > 2$) may be useful in dual problems.

6 Appendix

6.1 Derivation of the dual formulation (D) of (P1)

The derivation below is similar to [14, Lemma 2].

We replace $\max_{\mathbf{G}, \boldsymbol{\delta}} LR^2 \delta_{N-1}$ of (P1) by $\min_{\mathbf{G}, \boldsymbol{\delta}} \{-\delta_{N-1}\}$ for convenience in this section. The corresponding dual function of such (P1) is then defined as

$$H(\boldsymbol{\lambda}, \boldsymbol{\tau}; \mathbf{h}) = \min_{\substack{\mathbf{G} \in \mathbb{R}^{N \times d} \\ \boldsymbol{\delta} \in \mathbb{R}^N}} \{\mathcal{L}(\mathbf{G}, \boldsymbol{\delta}, \boldsymbol{\lambda}, \boldsymbol{\tau}; \mathbf{h}) := \mathcal{L}_1(\boldsymbol{\delta}, \boldsymbol{\lambda}, \boldsymbol{\tau}) + \mathcal{L}_2(\mathbf{G}, \boldsymbol{\lambda}, \boldsymbol{\tau}; \mathbf{h})\}$$

for dual variables $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_{N-1}]^\top \in \mathbb{R}_+^{N-1}$ and $\boldsymbol{\tau} = [\tau_0, \dots, \tau_{N-1}]^\top \in \mathbb{R}_+^N$, where $\mathcal{L}(\mathbf{G}, \boldsymbol{\delta}, \boldsymbol{\lambda}, \boldsymbol{\tau}; \mathbf{h})$ is a Lagrangian function, and

$$\begin{aligned} \mathcal{L}_1(\boldsymbol{\delta}, \boldsymbol{\lambda}, \boldsymbol{\tau}) &:= -\delta_{N-1} + \sum_{i=1}^{N-1} \lambda_i (\delta_i - \delta_{i-1}) + \sum_{i=0}^{N-1} \tau_i \delta_i, \\ \mathcal{L}_2(\mathbf{G}, \boldsymbol{\lambda}, \boldsymbol{\tau}; \mathbf{h}) &:= \sum_{i=1}^{N-1} \lambda_i \text{Tr}\{\mathbf{G}^\top \check{\mathbf{A}}_{i-1,i}(\mathbf{h}) \mathbf{G}\} + \sum_{i=0}^{N-1} \tau_i \text{Tr}\{\mathbf{G}^\top \check{\mathbf{D}}_i(\mathbf{h}) \mathbf{G} + \nu \mathbf{u}_i^\top \mathbf{G}\}. \end{aligned}$$

Here, $\min_{\boldsymbol{\delta}} \mathcal{L}_1(\boldsymbol{\delta}, \boldsymbol{\lambda}, \boldsymbol{\tau}) = 0$ for any $(\boldsymbol{\lambda}, \boldsymbol{\tau}) \in \Lambda$ where Λ is defined in (23), and $\min_{\boldsymbol{\delta}} \mathcal{L}_1(\boldsymbol{\delta}, \boldsymbol{\lambda}, \boldsymbol{\tau}) = -\infty$ otherwise.

For any given unit vector $\boldsymbol{\nu}$, [14, Lemma 1] implies

$$\min_{\mathbf{G} \in \mathbb{R}^{N \times d}} \mathcal{L}_2(\mathbf{G}, \boldsymbol{\lambda}, \boldsymbol{\tau}) = \min_{\mathbf{w} \in \mathbb{R}^N} \mathcal{L}_2(\mathbf{w} \boldsymbol{\nu}^\top, \boldsymbol{\lambda}, \boldsymbol{\tau}),$$

and thus for any $(\boldsymbol{\lambda}, \boldsymbol{\tau}) \in \Lambda$, we can rewrite the dual function as

$$\begin{aligned} H(\boldsymbol{\lambda}, \boldsymbol{\tau}; \mathbf{h}) &= \min_{\mathbf{w} \in \mathbb{R}^N} \{\mathbf{w}^\top \mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) \mathbf{w} + \boldsymbol{\tau}^\top \mathbf{w}\} \\ &= \max_{\gamma \in \mathbb{R}} \left\{ -\frac{1}{2} \gamma : \mathbf{w}^\top \mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) \mathbf{w} + \boldsymbol{\tau}^\top \mathbf{w} \geq -\frac{1}{2} \gamma, \forall \mathbf{w} \in \mathbb{R}^N \right\} \\ &= \max_{\gamma \in \mathbb{R}} \left\{ -\frac{1}{2} \gamma : \begin{pmatrix} \mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau}) & \frac{1}{2} \boldsymbol{\tau} \\ \frac{1}{2} \boldsymbol{\tau}^\top & \frac{1}{2} \gamma \end{pmatrix} \succeq 0 \right\}, \end{aligned}$$

where $\mathbf{S}(\mathbf{h}, \boldsymbol{\lambda}, \boldsymbol{\tau})$ is defined in (24). Therefore the dual problem of (P1) becomes (D), recalling that we previously replaced $\max_{\mathbf{G}, \boldsymbol{\delta}} LR^2 \delta_{N-1}$ of (P1) by $\min_{\mathbf{G}, \boldsymbol{\delta}} \{-\delta_{N-1}\}$.

6.2 Proof of Prop. 1

The proof is similar to [20, Prop. 2, 3 and 4].

We first show that $\{h_{i+1,k}\}$ in (25) is equivalent to

$$h_{i+1,k} = \begin{cases} \frac{(T_i - t_i)t_{i+1}}{t_i T_{i+1}} h_{i,k} & i = 0, \dots, N-1, k = 0, \dots, i-2, \\ \frac{(T_i - t_i)t_{i+1}}{t_i T_{i+1}} (h_{i,i-1} - 1), & i = 0, \dots, N-1, k = i-1, \\ 1 + \frac{(t_i - 1)t_{i+1}}{T_{i+1}}, & i = 0, \dots, N-1, k = i, \end{cases} \quad (48)$$

We use the notation $h'_{i,k}$ for the coefficients (25) to distinguish from (48). It is obvious that $h'_{i+1,i} = h_{i+1,i}$, $i = 0, \dots, N-1$, and we clearly have

$$\begin{aligned} h'_{i+1,i-1} &= \frac{t_{i+1}}{T_{i+1}} (t_{i-1} - h'_{i,i-1}) = \frac{t_{i+1}}{T_{i+1}} \left(t_{i-1} - \left(1 + \frac{(t_{i-1} - 1)t_i}{T_i} \right) \right) \\ &= \frac{(t_{i-1} - 1)(T_i - t_i)t_{i+1}}{T_i T_{i+1}} = \frac{(T_i - t_i)t_{i+1}}{t_i T_{i+1}} (h_{i,i-1} - 1) = h_{i+1,i-1} \end{aligned}$$

We next use induction by assuming $h'_{i+1,k} = h_{i+1,k}$ for $i = 0, \dots, n-1$, $k = 0, \dots, i$. We then have

$$h'_{n+1,k} = \frac{t_{n+1}}{T_{n+1}} \left(t_k - \sum_{j=k+1}^n h'_{j,k} \right) = \frac{t_{n+1}}{T_{n+1}} \left(t_k - \sum_{j=k+1}^{n-1} h'_{j,k} - h'_{n,k} \right)$$

$$= \frac{t_{n+1}}{T_{n+1}} \left(\frac{T_n}{t_n} h'_{n,k} - h'_{n,k} \right) = \frac{(T_n - t_n)t_{n+1}}{t_n T_{n+1}} h_{n,k} = h_{n+1,k}$$

Next, using (48), we show that FSFOM with (25) is equivalent to the GFPGM. We use induction, and for clarity, we use the notation $\mathbf{y}'_0, \dots, \mathbf{y}'_N$ for FSFOM with (48). It is obvious that $\mathbf{y}'_0 = \mathbf{y}_0$, and we have

$$\begin{aligned} \mathbf{y}'_1 &= \mathbf{y}'_0 - \frac{1}{L} h_{1,0} \tilde{\nabla}_L F(\mathbf{y}'_0) = \mathbf{y}_0 - \frac{1}{L} \left(1 + \frac{(t_0 - 1)t_1}{T_1} \right) \tilde{\nabla}_L F(\mathbf{y}_0) \\ &= \mathbf{x}_1 + \frac{(T_0 - t_0)t_1}{t_0 T_1} (\mathbf{x}_1 - \mathbf{x}_0) + \frac{(t_0^2 - T_0)t_1}{t_0 T_1} (\mathbf{x}_1 - \mathbf{y}_0) = \mathbf{y}_1, \end{aligned}$$

since $T_0 = t_0$. Assuming $\mathbf{y}'_i = \mathbf{y}_i$ for $i = 0, \dots, n$, we then have

$$\begin{aligned} \mathbf{y}'_{n+1} &= \mathbf{y}'_n - \frac{1}{L} h_{n+1,n} \tilde{\nabla}_L F(\mathbf{y}'_n) - \frac{1}{L} h_{n+1,n-1} \tilde{\nabla}_L F(\mathbf{y}'_{n-1}) - \frac{1}{L} \sum_{k=0}^{n-2} h_{n+1,k} \tilde{\nabla}_L F(\mathbf{y}'_k) \\ &= \mathbf{y}_n - \frac{1}{L} \left(1 + \frac{(t_n - 1)t_{n+1}}{T_{n+1}} \right) \tilde{\nabla}_L F(\mathbf{y}_n) \\ &\quad - \frac{1}{L} \frac{(T_n - t_n)t_{n+1}}{t_n T_{n+1}} (h_{n,n-1} - 1) \tilde{\nabla}_L F(\mathbf{y}_{n-1}) - \frac{1}{L} \sum_{k=0}^{n-2} \frac{(T_n - t_n)t_{n+1}}{t_n T_{n+1}} h_{n,k} \tilde{\nabla}_L F(\mathbf{y}_k) \\ &= \mathbf{x}_{n+1} - \frac{1}{L} \frac{(t_n^2 - T_n)t_{n+1}}{t_n T_{n+1}} \tilde{\nabla}_L F(\mathbf{y}_n) \\ &\quad - \frac{1}{L} \frac{(T_n - t_n)t_{n+1}}{t_n T_{n+1}} \left(\tilde{\nabla}_L F(\mathbf{y}_n) - \tilde{\nabla}_L F(\mathbf{y}_{n-1}) + \sum_{k=0}^{n-1} h_{n,k} \tilde{\nabla}_L F(\mathbf{y}_k) \right) \\ &= \mathbf{x}_{n+1} + \frac{(t_n^2 - T_n)t_{n+1}}{t_n T_{n+1}} (\mathbf{x}_{n+1} - \mathbf{y}_n) \\ &\quad + \frac{(T_n - t_n)t_{n+1}}{t_n T_{n+1}} \left(-\frac{1}{L} \tilde{\nabla}_L F(\mathbf{y}_n) + \frac{1}{L} \tilde{\nabla}_L F(\mathbf{y}_{n-1}) + \mathbf{y}_n - \mathbf{y}_{n-1} \right) \\ &= \mathbf{x}_{n+1} + \frac{(T_n - t_n)t_{n+1}}{t_n T_{n+1}} (\mathbf{x}_{n+1} - \mathbf{x}_n) + \frac{(t_n^2 - T_n)t_{n+1}}{t_n T_{n+1}} (\mathbf{x}_{n+1} - \mathbf{y}_n) = \mathbf{y}_{n+1}. \end{aligned}$$

6.3 Proof of Prop. 2

The proof is similar to [20, Prop. 1 and 5].

We use induction, and for clarity, we use the notation $\mathbf{y}'_0, \dots, \mathbf{y}'_N$ for FSFOM with (25) that is equivalent to GFPGM by Prop. 1. It is obvious that $\mathbf{y}'_0 = \mathbf{y}_0$, and we have

$$\begin{aligned} \mathbf{y}'_1 &= \mathbf{y}'_0 - \frac{1}{L} h_{1,0} \tilde{\nabla}_L F(\mathbf{y}'_0) = \mathbf{y}_0 - \frac{1}{L} \left(1 + \frac{(t_0 - 1)t_1}{T_1} \right) \tilde{\nabla}_L F(\mathbf{y}_0) \\ &= \left(1 - \frac{t_1}{T_1} \right) \left(\mathbf{y}_0 - \frac{1}{L} \tilde{\nabla}_L F(\mathbf{y}_0) \right) + \frac{t_1}{T_1} \left(\mathbf{y}_0 - \frac{1}{L} t_0 \tilde{\nabla}_L F(\mathbf{y}_0) \right) \\ &= \left(1 - \frac{t_1}{T_1} \right) \mathbf{x}_1 + \frac{t_1}{T_1} \mathbf{z}_1 = \mathbf{y}_1. \end{aligned}$$

Assuming $\mathbf{y}'_i = \mathbf{y}_i$ for $i = 0, \dots, n$, we then have

$$\begin{aligned} \mathbf{y}'_{n+1} &= \mathbf{y}'_n - \frac{1}{L} h_{n+1,n} \tilde{\nabla}_L F(\mathbf{y}'_n) - \frac{1}{L} \sum_{k=0}^{n-1} h_{n+1,k} \tilde{\nabla}_L F(\mathbf{y}'_k) \\ &= \mathbf{y}_n - \frac{1}{L} \left(1 + \frac{(t_n - 1)t_{n+1}}{T_{n+1}} \right) \tilde{\nabla}_L F(\mathbf{y}_n) - \frac{1}{L} \sum_{k=0}^{n-1} \frac{t_{n+1}}{T_{n+1}} \left(t_k - \sum_{j=k+1}^n h_{j,k} \right) \tilde{\nabla}_L F(\mathbf{y}_k) \end{aligned}$$

$$\begin{aligned}
&= \left(1 - \frac{t_{n+1}}{T_{n+1}}\right) \left(\mathbf{y}_n - \frac{1}{L} \tilde{\nabla}_L F(\mathbf{y}_n)\right) \\
&\quad + \frac{t_{n+1}}{T_{n+1}} \left(\mathbf{y}_n - \frac{1}{L} \sum_{k=0}^n t_k \tilde{\nabla}_L F(\mathbf{y}_k) + \frac{1}{L} \sum_{k=0}^{n-1} \sum_{j=k+1}^n h_{j,k} \tilde{\nabla}_L F(\mathbf{y}_k)\right) \\
&= \left(1 - \frac{t_{n+1}}{T_{n+1}}\right) \left(\mathbf{y}_n - \frac{1}{L} \tilde{\nabla}_L F(\mathbf{y}_n)\right) + \frac{t_{n+1}}{T_{n+1}} \left(\mathbf{y}_0 - \frac{1}{L} \sum_{k=0}^n t_k \tilde{\nabla}_L F(\mathbf{y}_k)\right) \\
&= \left(1 - \frac{t_{n+1}}{T_{n+1}}\right) \mathbf{x}_{n+1} + \frac{t_{n+1}}{T_{n+1}} \mathbf{z}_{n+1}.
\end{aligned}$$

6.4 Appendix of footnote 7

Our formulation (P') examines the set $\{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}, \mathbf{x}_N\}$ and eventually leads to the best known analytical bound on the norm of the composite gradient mapping in Thm. 5 among fixed-step first-order methods. An alternative formulation would be to use the set $\{\mathbf{y}_0, \dots, \mathbf{y}_{N-1}\}$ (i.e., excluding the point \mathbf{x}_N). For this alternative, we could simply replace the inequality (33) with the condition $0 \leq F(\mathbf{y}_{N-1}) - F(\mathbf{x}_*)$ to derive a slightly different relaxation. (One could use other conditions using a subgradient at the point \mathbf{y}_{N-1} , but this is beyond the scope of this paper.) However, we found that this approach led to a larger upper bound than (40). Another alternative would be to use the set $\{\mathbf{x}_0, \dots, \mathbf{x}_N\}$, but deriving a tight relaxation for this case remains an open problem. Nevertheless, the inequality in Lemma 1 provides a bound for that set as seen in Thm. 4 and 5. Another alternative would be simply to use the final point $\{\mathbf{x}_N\}$ or $\{\mathbf{y}_N\}$ in (P') instead of the minimum over a set of points. However, the corresponding relaxation (P1') yielded only an $O(1/N)$ bound at best on the final composite gradient mapping norm in this case, which is worse than that of FPGM- m in (15). So we leave finding its tighter relaxation as future work.

6.5 Proof of Equation (45) in Thm. 5

$$\begin{aligned}
&\sum_{k=0}^{N-1} (T_k - t_k^2) + T_{N-1} \\
&= \sum_{k=m}^{N-1} \left(t_{m-1}^2 + \sum_{l=m}^k t_l - t_k^2\right) + t_{m-1}^2 + \sum_{l=m}^{N-1} t_l \\
&= (N-m+1)t_{m-1}^2 + \sum_{k=m}^{N-1} \left(\sum_{l=m}^k \frac{N-l+1}{2} - \left(\frac{N-k+1}{2}\right)^2\right) + \sum_{l=m}^{N-1} \frac{N-l+1}{2} \\
&= (N-m+1)t_{m-1}^2 + \sum_{k'=0}^{N-m-1} \left(\sum_{l'=0}^{k'} \frac{N-l'-m+1}{2} - \left(\frac{N-k'-m+1}{2}\right)^2\right) \\
&\quad + \sum_{l'=0}^{N-m-1} \frac{N-l'-m+1}{2} \\
&= (N-m+1)t_{m-1}^2 \\
&\quad + \sum_{k=0}^{N-m-1} \left(\frac{(N-m+1)(k+1)}{2} - \frac{k(k+1)}{4} - \frac{(N-m+1)^2 - 2(N-m+1)k + k^2}{4}\right) \\
&\quad + \frac{(N-m+1)(N-m)}{2} - \frac{(N-m-1)(N-m)}{4} \\
&= (N-m+1)t_{m-1}^2 + \sum_{k=0}^{N-m-1} \left(-\frac{k^2}{2} + (N-m+3/4)k - \frac{(N-m-1)(N-m+1)}{4}\right)
\end{aligned}$$

$$\begin{aligned}
& + \frac{(N-m)(N-m+3)}{4} \\
& = (N-m+1)t_{m-1}^2 - \frac{(N-m-1)(N-m-1/2)(N-m)}{6} \\
& \quad + \frac{(N-m-1)(N-m)(N-m+3/4)}{2} \\
& \quad - \frac{(N-m-1)(N-m)(N-m+1)}{4} + \frac{(N-m)(N-m+3)}{4} \\
& \geq \frac{(N-m+1)(m+1)^2}{4} + \frac{(N-m-1)(N-m)^2}{3} - \frac{(N-m)^2(N-m+1)}{4} \\
& \geq \frac{(N-m-1)(N-m)^2}{3} \\
& \geq \frac{1}{24}(N-2)N^2,
\end{aligned}$$

where $m = \lfloor \frac{N}{2} \rfloor \geq \frac{N-1}{2}$, $N-m \geq \frac{N}{2}$, and $t_{m-1} \geq \frac{m+1}{2}$ (13).

Acknowledgements

The authors would like to thank the anonymous referees for very useful comments that have improved the quality of this paper.

References

1. H. ATTOUCH, Z. CHBANI, J. PEYPOUQUET, AND P. REDONT, *Fast convergence of inertial dynamics and algorithms with asymptotic vanishing viscosity*, Mathematical Programming, (2016), doi:10.1007/s10107-016-0992-8.
2. H. ATTOUCH AND J. PEYPOUQUET, *The rate of convergence of Nesterov's accelerated forward-backward method is actually faster than $1/k^2$* , SIAM J. Optim., 26 (2016), pp. 1824–34, doi:10.1137/15M1046095.
3. A. BECK, A. NEDIC, A. OZDAGLAR, AND M. TEBoulLE, *An $O(1/k)$ gradient method for network resource allocation problems*, IEEE Trans. Control of Network Systems, 1 (2014), pp. 64–73, doi:10.1109/TCNS.2014.2309751.
4. A. BECK AND M. TEBoulLE, *Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems*, IEEE Trans. Im. Proc., 18 (2009), pp. 2419–34, doi:10.1109/TIP.2009.2028250.
5. A. BECK AND M. TEBoulLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM J. Imaging Sci., 2 (2009), pp. 183–202, doi:10.1137/080716542.
6. A. BECK AND M. TEBoulLE, *A fast dual proximal gradient algorithm for convex minimization and applications*, Operations Research Letters, 42 (2014), pp. 1–6, doi:10.1016/j.orl.2013.10.007.
7. S. BUBECK, Y. T. LEE, AND M. SINGH, *A geometric alternative to nesterov's accelerated gradient descent*, 2015, <http://arxiv.org/abs/1506.08187>. arxiv 1506.08187.
8. A. CHAMBOLLE AND C. DOSSAL, *On the convergence of the iterates of the "Fast iterative shrinkage/Thresholding algorithm"*, J. Optim. Theory Appl., 166 (2015), pp. 968–82, doi:10.1007/s10957-015-0746-4.
9. P. L. COMBETTES AND J.-C. PESQUET, *Proximal splitting methods in signal processing*, 2011, doi:10.1007/978-1-4419-9569-8_10. Fixed-Point Algorithms for Inverse Problems in Science and Engineering, Springer, Optimization and Its Applications, pp 185–212.
10. I. CVX RESEARCH, *CVX: Matlab software for disciplined convex programming, version 2.0*. <http://cvxr.com/cvx>, Aug. 2012.
11. I. DAUBECHIES, M. DEFRISE, AND C. DE MOL, *An iterative thresholding algorithm for linear inverse problems with a sparsity constraint*, Comm. Pure Appl. Math., 57 (2004), pp. 1413–57, doi:10.1002/cpa.20042.
12. O. DEVOLDER, F. GLINEUR, AND Y. NESTEROV, *Double smoothing technique for large-scale linearly constrained convex optimization*, SIAM J. Optim., 22 (2012), pp. 702–27, doi:10.1137/110826102.
13. Y. DRORI, *The exact information-based complexity of smooth convex minimization*, Journal of Complexity, 39 (2017), pp. 1–16, doi:10.1016/j.jco.2016.11.001.
14. Y. DRORI AND M. TEBoulLE, *Performance of first-order methods for smooth convex minimization: A novel approach*, Mathematical Programming, 145 (2014), pp. 451–82, doi:10.1007/s10107-013-0653-0.
15. Y. DRORI AND M. TEBoulLE, *An optimal variant of Kelley's cutting-plane method*, Mathematical Programming, 160 (2016), pp. 321–51, doi:10.1007/s10107-016-0985-7.
16. S. GHADIMI AND G. LAN, *Accelerated gradient methods for nonconvex nonlinear and stochastic programming*, Mathematical Programming, 156 (2016), pp. 59–99, doi:10.1007/s10107-015-0871-8.
17. T. GOLDSTEIN, B. O'DONOGHUE, S. SETZER, AND R. BARANIUK, *Fast alternating direction optimization methods*, SIAM J. Imaging Sci., 7 (2014), pp. 1588–623, doi:10.1137/120896219.
18. M. GRANT AND S. BOYD, *Graph implementations for nonsmooth convex programs*, in Recent Advances in Learning and Control, V. Blondel, S. Boyd, and H. Kimura, eds., Lecture Notes in Control and Information Sciences, Springer-Verlag Limited, 2008, pp. 95–110. http://stanford.edu/~boyd/graph_dcp.html.

19. D. KIM AND J. A. FESSLER, *Generalizing the optimized gradient method for smooth convex minimization*, 2016, <http://arxiv.org/abs/1607.06764>. arxiv 1607.06764.
20. D. KIM AND J. A. FESSLER, *Optimized first-order methods for smooth convex minimization*, Mathematical Programming, 159 (2016), pp. 81–107, [doi:10.1007/s10107-015-0949-3](https://doi.org/10.1007/s10107-015-0949-3).
21. D. KIM AND J. A. FESSLER, *On the convergence analysis of the optimized gradient method*, J. Optim. Theory Appl., 172 (2017), pp. 187–205, [doi:10.1007/s10957-016-1018-7](https://doi.org/10.1007/s10957-016-1018-7).
22. L. LESSARD, B. RECHT, AND A. PACKARD, *Analysis and design of optimization algorithms via integral quadratic constraints*, SIAM J. Optim., 26 (2016), pp. 57–95, [doi:10.1137/15M1009597](https://doi.org/10.1137/15M1009597).
23. R. D. C. MONTEIRO AND B. F. SVAITER, *An accelerated hybrid proximal extragradient method for convex optimization and its implications to second-order methods*, SIAM J. Optim., 23 (2013), pp. 1092–1125, [doi:10.1137/110833786](https://doi.org/10.1137/110833786).
24. I. NECOARA AND A. PATRASCU, *Iteration complexity analysis of dual first order methods for conic convex programming*, Optimization Methods and Software, 31 (2016), pp. 645–78, [doi:10.1080/10556788.2016.1161763](https://doi.org/10.1080/10556788.2016.1161763).
25. Y. NESTEROV, *A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$* , Dokl. Akad. Nauk. USSR, 269 (1983), pp. 543–7.
26. Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, Kluwer, 2004, <http://books.google.com/books?id=VyYLem-13CgC>.
27. Y. NESTEROV, *Smooth minimization of non-smooth functions*, Mathematical Programming, 103 (2005), pp. 127–52, [doi:10.1007/s10107-004-0552-5](https://doi.org/10.1007/s10107-004-0552-5).
28. Y. NESTEROV, *How to make the gradients small*, Optima, 88 (2012), pp. 10–11.
29. Y. NESTEROV, *Gradient methods for minimizing composite functions*, Mathematical Programming, 140 (2013), pp. 125–61, [doi:10.1007/s10107-012-0629-5](https://doi.org/10.1007/s10107-012-0629-5).
30. B. O'DONOGHUE AND E. CANDÈS, *Adaptive restart for accelerated gradient schemes*, Found. Comp. Math., 15 (2015), pp. 715–32, [doi:10.1007/s10208-013-9150-3](https://doi.org/10.1007/s10208-013-9150-3).
31. W. SU, S. BOYD, AND E. J. CANDÈS, *A differential equation for modeling Nesterov's accelerated gradient method: theory and insights*, J. Mach. Learning Res., 17 (2016), pp. 1–43.
32. A. B. TAYLOR, J. M. HENDRICKX, AND F. GLINEUR, *Exact worst-case performance of first-order algorithms for composite convex optimization*, 2015, <http://arxiv.org/abs/1512.07516>. arxiv 1512.07516.
33. A. B. TAYLOR, J. M. HENDRICKX, AND F. GLINEUR, *Smooth strongly convex interpolation and exact worst-case performance of first-order methods*, Mathematical Programming, 161 (2017), pp. 307–45, [doi:10.1007/s10107-016-1009-3](https://doi.org/10.1007/s10107-016-1009-3).